

(19) 日本国特許庁 (J P)

## (12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平9-200239

(43) 公開日 平成9年(1997)7月31日

(51) Int.Cl. <sup>8</sup>	識別記号	庁内整理番号	F I	技術表示箇所
H 0 4 L 12/42			H 0 4 L 11/00	3 3 0
G 0 6 F 13/00	3 5 7		G 0 6 F 13/00	3 5 7 C
13/36	5 3 0		13/36	5 3 0 C

審査請求 未請求 請求項の数9 O L (全 44 頁)

(21) 出願番号 特願平8-7140

(22) 出願日 平成8年(1996)1月19日

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72) 発明者 岡田 康行

神奈川県海老名市下今泉810番地株式会社

日立製作所オフィスシステム事業部内

(74) 代理人 弁理士 小川 勝男

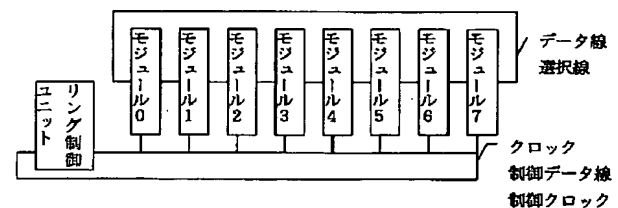
(54) 【発明の名称】 リング接続を用いたデータ転送方法及び情報処理システム

(57) 【要約】

【課題】長距離のプロセッサ間の接続方法であったリング接続をプロセッサ間的高速接続に適用する際の隘路となっていたレイテンシ、スループット、運用機能を改善し、プロセッサ間を高速に接続する。

【解決手段】複数の信号線からなるリングでプロセッサ等のモジュールを複数接続し、送信権を獲得するためのフラグ用の独立する信号線を設け、前記フラグをセットして送信権を要求し、受信したフラグがセットされていないことから前記送信権を先行して獲得したことを後から確認することで受信及び送信動作と前記リングの転送動作を並行して行う。

図24



## 【特許請求の範囲】

【請求項 1】複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおけるリング接続を用いたデータ転送方法であって、送信権を獲得するためのフラグ用の独立する信号線を設け、前記フラグをセットして送信権を要求し、受信したフラグがセットされていないことから前記送信権を先行して獲得したことを後から確認することで受信及び送信動作と前記リングの転送動作を並行して行うことを特徴とするリング接続を用いたデータ転送方法。

【請求項 2】複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおけるリング接続を用いたデータ転送方法であって、同一クロックを前記リングに配信し、転送及び送受信に使用することを特徴とするリング接続を用いたデータ転送方法。

【請求項 3】複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおけるリング接続を用いたデータ転送方法であって、位相調整回路によりクロックの位相を受信する信号に合わせ、前記クロックの乗り換えの遅延を不要にするためクロックを転送制御、受信制御、送信制御に同一位相で使用し、前記複数のモジュールのうちの少なくとも一つのモジュールは回送モードで前記リング 1 周の遅延を前記クロックの整数倍のサイクルに合わせることを特徴とするリング接続を用いたデータ転送方法。

【請求項 4】複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおけるリング接続を用いたデータ転送方法であって、前記複数のモジュールのうちの各モジュールはラッチせずに前記複数のモジュールのうちの次のモジュールに転送し、ラッチ前の位相に合わせて送信することを特徴とするリング接続を用いたデータ転送方法。

【請求項 5】複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおけるリング接続を用いた情報処理システムであって、送信権を獲得するためのフラグ用の独立する信号線を設け、前記フラグをセットして送信権を要求し、受信したフラグがセットされていないことから前記送信権を先行して獲得したことを後から確認する手段を有することを特徴とするリング接続を用いた情報処理システム。

【請求項 6】複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおけるリング接続を用いた情報処理システムであって、クロックの位相を受信する信号に合わせる位相調整回路と、転送制御、受信制御、送信制御に同一位相で使用するクロックを供給する手段と、前記リング 1 周の遅延を前記クロックの整数倍のサイクルに合わせる少なくとも一つの回送モードのモジュールを有することを特徴とするリング接続を用いた情報処理システム。

【請求項 7】複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおけるリング接続を用いた情報処理システムであって、前記複数のモジュールのうちの各モジュールは、ラッチせずに前記複数のモジュールのうちの次のモジュールに転送する手段と、ラッチ前の位相に合わせて送信する手段とを有することを特徴とするリング接続を用いた情報処理システム。

10 【請求項 8】複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおけるリング接続を用いた情報処理システムであって、前記複数のモジュールのうちの各モジュールは、ラッチせずに前記複数のモジュールのうちの次のモジュールに転送する手段と、ラッチ前の位相に合わせて送信する手段とを有することを特徴とするリング接続を用いた情報処理システム。

20 【請求項 9】複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおけるリング接続を用いた情報処理システムであって、パケットを連続的に転送するブロック・パケットの送信権を獲得するための信号線と、前記送信権を獲得した 2 サイクル後からブロック転送の間通常のパケットの送信要求を抑止させことを通知する信号線と、ブロックパケットの送信権の獲得後の 4 サイクル後からブロックパケットの送信を許可する手段とを有することを特徴とするリング接続を用いた情報処理システム。

## 【発明の詳細な説明】

## 【0001】

30 【発明の属する技術分野】本発明は、パーソナルコンピュータ、ワークステーション、サーバ等の情報処理装置に用いられるデータ転送方法及びそれを用いた情報処理システムに関し、特に複数のプロセッサ間のデータ転送を高速に行うのに好適なリング接続を用いたデータ転送方法及びそれを用いた情報処理システムに関する。

## 【0002】

40 【従来の技術】パーソナルコンピュータ等の情報処理装置が普及し、様々な分野で使用されるようになるにつれて、情報処理装置の処理性能の向上が要求されている。性能向上の手段として、情報処理装置を複数のプロセッサから構成するアプローチがある。この構成の場合、複数のプロセッサ間のデータ転送を高速に行うことが性能上の重要な鍵となる。プロセッサ間の高速接続手段として、リング接続方式と現在使用されている接続方式であるバス、スイッチとの比較を図 1 に示す。

50 【0003】プロセッサ間データ転送を高速に行うためのプロセッサ間の接続手段として、バスが最も広く用いられている。バスは外部に付加回路を使用せず接続でき、拡張性に富むとともに、バス制御 LSI、コネクタなどに必要なピンの数が少なく低価格である。また、すべてのプロセッサは、バスの内容を見ることができた

め密結合のプロセッサ間接続の全報知の通信（スヌープ）に適する。

【0004】バス接続の場合、性能、接続台数を決める主要な要因である通信容量は、動作周波数に依存する。プロセッサバスでは直近のプロセッサと最遠のプロセッサからの信号伝搬の時間差を同一クロック内に納める必要があり、動作周波数にはバスの長さによる限界がある。クロック、LSI、伝送線のスキューの設計値は、動作周波数に直接影響を与える。送信元と受信先にすべての組み合わせがあるため、スキュー補正に適切な手段がない。プロセッサへの引出線からの信号の両方向の伝搬、引出線の通過により伝搬特性の劣化、すべてのプロセッサへの駆動能力が必要なことによる動作周波数の限界もある。バスの長さを短くして動作周波数を向上させることも行われているが、マイクロプロセッサ・システムの実装形態に制約を与えてしまう。プロセッサのクロックとデータを一緒に送信すれば動作周波数は向上できるが、プロセッサバスで頻発するプロセッサ間の切り替え時間は増加する。接続できるプロセッサ数は電氣的に限界がある。また、バスの送信権を獲得するための遅れは別の性能限界を与える。さらに、プロセッサ間の転送の切り替えが頻繁で、プロセッサからの応答が一定しないプロセッサ間の高性能の接続には課題が残されている。

【0005】プロセッサ間をスイッチで接続する場合、プロセッサとスイッチとの間で1方向に転送し、1対1の通信であることから動作周波数を向上できる。全体の通信容量は方式的には限界がなく、接続台数の制約はない。一方で、スイッチは外部に集中させる必要があるため接続はスター状になり、総配線量が多く、スイッチ周辺の配線が込み合う。また、全報知の通信は制御を複雑にして性能を劣化させる。スイッチを経由する遅延時間も性能を制約する。

【0006】プロセッサ間の接続には、リングも用いられる。リング接続は、信号を一方に1対1に転送すればよいので、バスより通信容量を大きくとれる。次のプロセッサまでの伝送線の長さが短く、クロック、LSI、伝送線のスキューの補正を可能にし、信号の引出線が不要であり、次のプロセッサだけを駆動すればよいので動作周波数の制約を軽減できる可能性がある。また、バスのような接続の長さの制約は少なく接続台数の電氣的制約もない。バスと同じく付加回路は必要なく、バスと同様に接続の配線量は少ない。さらに、すべてのプロセッサはリングの内容を見ることができるため、密結合のプロセッサ間接続の全報知の通信（スヌープ）にも適する。

【0007】一方で、リング接続の場合、バスに対して2倍のピン数を必要とする。バスでは1クロックですべてのプロセッサに送信できるが、リングではプロセッサを経由する時間遅れがスキュー補正、クロック同期、送

信権の制御等のため大きく、特に密結合のプロセッサ間接続には課題がある。

【0008】プロセッサ間に応用したリングとして、例えば、IEEE Standard 1596-1992, "Scalable Coherent Interface (SCI)," 1992、あるいは、D. Cecchi, M. Dina, C. Preuss, "A 1GB/S SCI Link in 0.8μMBi CMOS," 1995 IEEE International Solid-State Circuits Conference, San Francisco, CA, Digest of Technical Papers. Paper 20.2, February, 1995, pp. 326-327. に記載のものがある。

【0009】上記文献に記載されている代表的なプロセッサ間接続であるSCIは、上述の課題を十分には解決していない。

【0010】プロセッサ間接続手段としてリング接続を用いた場合のリングのインターフェース部の構成を図2に示す。

【0011】図2において、バッファ801から受信して、フラグ及びアドレスを判定して受信FIFO807、或いは中継用FIFO806に送る処理をする。受信クロックで動作する範囲802とプロセッサ等の当該モジュールのクロックは違うものを使用するので速度差吸収用のバッファである受信FIFO807、及び中継用FIFO806が必要である。受信クロックと当該モジュールのクロックは当然位相も異なるのでバッファ801と送受処理805の間にクロックの乗り換えが必要である。送信時には転送ラッチ803を経由して送信する。なお、時々々の同期パターンによってバッファ801に入る前に一括してスキュー補正される。送達確認の packets は、自動的に生成されて送信される。リングは1つであり、すべての情報交換は同一の packets 形式によって行われる。スヌープなど全報知の機能を用いずにキャッシュコヒーレンスを行う。リングの障害検出、切り離し、復旧は、通常の packets の転送パスを通じて行う。

【0012】

【発明が解決しようとする課題】従来技術に述べた通り、プロセッサバスに適合するためにはレイテンシ、スループット、運用性を改善し、またプロセッサバス固有の機能を効率よく実現する必要がある。そして、リング接続の場合のプロセッサを経由するための時間遅れを改善し、比較的近距离の接続を前提にするプロセッサ・バスの性能限界を超えることを可能とする必要がある。

【0013】本発明の目的は、プロセッサを経由するためのレイテンシを短縮することを可能とするリング接続を用いたデータ転送方法及び情報処理システムを提供することにある。

## 5

【0014】本発明の他の目的は、リングが通常動作しなくなる切り替え時にもリングよりモジュールに安定したクロックを供給することを可能とするリング接続を用いた情報処理システムを提供することにある。

【0015】本発明の他の目的は、情報処理システムのスループットを向上させるリング接続を用いたデータ転送方法及び情報処理システムを提供することにある。

【0016】本発明の他の目的は、リングの構成変更に伴うシステム停止時間を短縮するリング接続を用いたデータ転送方法及び情報処理システムを提供することにある。

【0017】

【課題を解決するための手段】

【プロセッサバスへの適用】上記目的を達成するため、本発明のリング接続を用いたデータ転送方法は、複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおいて、送信権を獲得するためのフラグ用の独立する信号線を設け、前記フラグをセットして送信権を要求し、受信したフラグがセットされていないことから前記送信権を先行して獲得したことを後から確認することで受信及び送信動作と前記リングの転送動作を並行して行う。

【0018】別の観点からは、本発明のリング接続を用いたデータ転送方法は、複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおいて、同一クロックを前記リングに配信し、転送及び送受信に使用する。

【0019】さらに別の観点からは、本発明のリング接続を用いたデータ転送方法は、複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおいて、位相調整回路によりクロックの位相を受信する信号に合わせ、クロックの乗り換えの遅延を不要にするためクロックを転送制御、受信制御、送信制御に同一位相で使用し、前記複数のモジュールのうちの少なくとも一つのモジュールは回送モードで前記リング1周の遅延を前記クロックの整数倍のサイクルに合わせる。

【0020】さらに別の観点からは、本発明のリング接続を用いたデータ転送方法は、複数の信号線からなるリングでプロセッサ等のモジュールを複数接続した情報処理システムにおいて、複数のモジュールのうちの各モジュールはラッチせずに次のモジュールに転送し、ラッチ前の位相に合わせて送信する。

【0021】さらに別の観点からは、本発明のリング接続を用いた情報処理システムは、複数の信号線からなるリングでプロセッサ等のモジュールを複数接続し、送信権を獲得するためのフラグ用の独立する信号線を設け、前記フラグをセットして送信権を要求し、受信したフラグがセットされていないことから前記送信権を先行して獲得したことを後から確認する手段を有する。

## 6

【0022】さらに別の観点からは、本発明のリング接続を用いた情報処理システムは、複数の信号線からなるリングでプロセッサ等のモジュールを複数接続し、クロックの位相を受信する信号に合わせる位相調整回路と、転送制御、受信制御、送信制御に同一位相で使用するクロックを供給する手段と、前記リング1周の遅延を前記クロックの整数倍のサイクルに合わせる少なくとも一つの回送モードのモジュールを有する。

【0023】さらに別の観点からは、本発明のリング接続を用いた情報処理システムは、複数の信号線からなるリングでプロセッサ等のモジュールを複数接続し、複数のモジュールのうちの各モジュールは、ラッチせずに複数のモジュールのうちの次のモジュールに転送する手段と、ラッチ前の位相に合わせて送信する手段とを有する。

【0024】さらに別の観点からは、本発明のリング接続を用いた情報処理システムは、複数の信号線からなるリングでプロセッサ等のモジュールを複数接続し、複数のモジュールのうちの各モジュールは、ラッチせずに前記複数のモジュールのうちの次のモジュールに転送する手段と、ラッチ前の位相に合わせて送信する手段とを有する。

【0025】さらに別の観点からは、本発明のリング接続を用いた情報処理システムは、複数の信号線からなるリングでプロセッサ等のモジュールを複数接続し、パケットを連続的に転送するブロック・パケットの送信権を獲得するための信号線と、送信権を獲得した2サイクル後からブロック転送の間通常のパケットの送信要求を抑制させことを通知する信号線と、ブロックパケットの送信権の獲得後の4サイクル後からブロックパケットの送信を許可する手段とを有する。

【0026】

【発明の実施の形態】以下、本発明の実施例を図面を参照して詳細に説明する。

【0027】最初に、本発明の基本的な概念を説明する。

【0028】〔リングのプロセッサバスへの適用〕本発明においては、高速の遠距離の伝送手段として利用されているリングをその高速性を活かしてバスの性能限界を超えるためにプロセッサバスに適用する。プロセッサバスの最も重要な性能要素は遠距離の伝送では余り重要でないレイテンシの短縮である。一方、リングは転送の容量を共有するので本質的にスループットのネックとなる。また、遠距離の伝送では伝送系の障害回復が重要で、伝送系は独立している。それに対して、プロセッサバスでは接続モジュール数が少ない近距離の伝送であり、クロック、電源、冷却、筐体を共有する等から伝送系の障害は少ない。したがって、モジュールの障害検出、切り離し、構成変更等リングを含むシステムの一体運用が重要であり、レイテンシ、スループット、運用性

を改善してリングをプロセッサバスに適用する。

【0029】〔1サイクル1パケット転送〕1サイクル1パケット転送によりレイテンシの3サイクル短縮を図る。現状のリングは遠距離の伝送を前提にしており、転送のビット幅は距離の制約から1ビットから16ビット程度である。図3に示すように転送のビット幅が16ビットで1パケットが64ビットならば4サイクルになる。本発明のリングは近距離の伝送を前提にし、1サイクルで1パケットを1度に転送する方式をとる。図4に示すように1パケットを64ビットにすると転送のビット幅は64ビットで1サイクルになる。転送の方式に柔軟性を持たせるため、後述するようにブロックパケット、複数サイクルのパケットの転送も可能に拡張している。

【0030】〔送信権の先行獲得〕送信権の先行獲得によりレイテンシの1サイクル短縮を図る。現状のリングは図5に示すようにフラグを読みその内容に従い送信権を獲得してフラグをセットして次のモジュールにフラグを送信する。本発明のリングは、送信したいならばフラグをセットして図6に示す送信、論理和、転送ラッチを経由して次のモジュールにフラグを送信する。送信権を先行して獲得したことを受信したフラグがセットされていないことで後で確認する。

【0031】送信したモジュールはリング1周の時に合わせて、受信したフラグの転送を送信権を放棄との論理積により抑止してフラグをリセットする。各モジュールは図7に示すようにすべての情報を受け取りパケット内のヘッダを解読して自分宛のパケットだけを選択して受信し、並行してすべての情報を転送する。この方式により、読み、処理、書くという処理とパケットの転送を独立、並行させられる。パケットの送信は、あらかじめ送信権を獲得してその2サイクル後に行う。パケットを次のモジュールに単に転送すれば良いのでレイテンシを改善できる。

【0032】〔同一クロック配信〕同一クロック配信によりレイテンシの3-4サイクルの短縮を図る。現状のリングは、各モジュールが遠距離にあることを前提にしているので独立するクロックを使用しており、図8に示すようにクロックの速度差の吸収のためのバッファがいる。本発明では、制御リングによりクロックを配信することにより速度差の吸収バッファを不要にする。

【0033】〔受信位相での処理〕受信位相で処理することによりレイテンシの2分の1サイクル短縮を図る。現状のリングは、クロックの位相がモジュール毎に独立しているので、速度差吸収バッファに加え、図9に示すようにモジュール内で使用する時にクロックの乗り換えの遅延がいる。本発明では、図10に示すようにクロック位相調整回路により受信するタイミングを合わせ、これを処理、送信に使用するのでクロックの乗り換えが不要になる。

【0034】このようにすることにより、動作周波数はクロックのスキューに影響されなくなるので向上する。クロックの位相は、モジュールを経由する毎に遅延量に応じて進む。リング1周して再送信するときにはどのモジュールにおいてもクロックの整数倍になる必要がある。図11に示すように、リングに最低一つのモジュールは、回送モードを設けてリングの1周のサイクル数をクロックの整数倍（ここでは6サイクル）にするため、クロックに同期して送信する。本発明では、伝送系の変動があまりないことを前提にしており、動作中にテストパケットを送信して位相の調整を動的に行うのではなく、エラー発生時のみ再調整する方法を採用する。

【0035】〔位相変化時のクロックの切り替え〕位相変化時のクロックの切り替えによりクロックを連続供給する。モジュールの挿抜時、イニシャライズ時には、クロックの波形変動による影響を避け、連続的に動作を継続するため、PLLを経由してクロックを配信して、リング制御に関係する部分以外のモジュールの制御に用いモジュールの動作を保証する。図12に本発明のクロックシステムを示す。

【0036】〔ラッチせず中継する方式〕ラッチせず中継することによりレイテンシを2分の1サイクル短縮する。現状のリングでは図13に示すようにラッチを経由して受信パケットを転送する。本発明のリングでも、波形整形のため転送ラッチを経由して受信パケットを転送する。本発明では、図14に示すように、波形整形の必要がない場合は転送ラッチを経由せずに中継できる。ラッチは受信の立ち上がりの2分の1のタイミングで行うのでこれにより2分の1サイクルのレイテンシの短縮ができる。

【0037】送信パケットの位相は、転送ラッチで位相の調整ができないので、受信の立ち上がり位相を合わせるためのラッチを経由して送信する。本発明の各モジュールは、ラッチせず中継する中継モード、波形整形のため数モジュールに1つ必要になる転送ラッチを経由して中継する転送モード、リングに最低1つ必要なクロックに同期させる回送ラッチを経由する回送モードのいずれかのモードで動作する。図15にパケットA、B、Cを転送する例を示す。

【0038】〔受信線毎のスキュー補正〕受信線毎のスキュー補正により動作周波数を向上させる。現状のリングは、スキューを一括して補正してラッチしているため動作周波数はスキュー量が直接的に影響する。本発明では、簡単な回路により図16に示すように信号線毎にスキュー補正するので動作周波数は補正精度で決まる所まで向上できる。

【0039】〔信号線を用いた送達確認〕信号線を用いた送達確認によりスループットを向上させる。現状のリングは、送達確認を応答パケットによっている。本発明のリングは、応答パケットを送信する代わりに送達確認

のために送達確認とその報告を信号線で行うのでスループットを改善できる。送達確認をパケット毎に行う方式を取れば、スループットは2分の1になる。送信モジュールは、送信権獲得時に到達確認の報告と通知の2つの信号線の送信権を獲得する。受信モジュールは、パケットの受信の2サイクル後に送達確認を信号線で報告する。全報知についても送達確認を行うために受信できた時には何も応答しない応答方式を取るので、信号線1本で送達確認ができる。

【0040】[短いデータ幅のリングの併設] 短いデータ幅のリングを併設することによりスループットを向上させる。現状のリングは、遠距離の転送を目的としてきたため、用途別に短いデータ幅のリングを併設するという概念はない。短いパケットもリングを共用していた。プロセッサバスに適用するとデータ転送以外の短いパケットが多量に必要な場合がある。データ転送と共用するとスループット、転送効率が下がる。本発明のリングは、信号1本で送信権の制御ができるのでブロックパケットの転送を行うブロックパケットリング、1サイクルで2つのパケットを並列転送するための選択リング、キャッシュのディレクトリの読み出すためのディレクトリ・リング等短いデータ幅のパケットを持つリングを容易に併設できる。データ転送のためのリングを基本リングと呼ぶ。図17にリングの種類を示す。基本リングと併設リングはクロックを共通にすることにより、情報の交換、相互の動作の同期を行う。制御リングも共用するのでリング併設による複雑性を除去する。

【0041】[ブロックパケットリング] ブロックパケットリングによりスループットを向上させる。ブロックパケットリングは、パケットを連続的に転送して各モジュールの制御を簡単にし、またスループットを向上させるためにある。現状のリングでは、パケットのヘッダに連続転送のフラグに相当するものがある。この方式を取ると受信パケットを読み、処理、書くことで送信権を獲得し、通知することになり本方式で述べたレイテンシの短縮ができない。本リングの特徴を活かして、基本リングを拡張してブロックパケットを転送できる。

【0042】基本リングは、1サイクルに1パケットを転送している。まず、パケットの連続転送するモジュールを1つ選択する必要がある。選択のための送信権を獲得する機能を設ける。パケットの連続転送中は1サイクルに1パケットの転送を抑止しなければならない。これには、あらかじめパケットの連続送信を各モジュールに通知すればよい。各モジュールは、予約を通知された2サイクル後からパケットの送信要求を抑止する。かくして、送信モジュールは連続送信を通知した2サイクル後からパケットの送信要求を行うことができる。このため、図18に示すようにパケットの連続転送の送信権獲得のための信号である予約線、及び予約通知線の2本の信号からなるブロックパケットリングを併設すればパケ

ットの連続転送を実現できる。

【0043】[選択リング] 選択リングによりスループットを向上させる。基本リングは、送信権をリング1周の間獲得する。送信モジュールから受信モジュールまでパケットを転送する場合、受信モジュールから送信モジュールまで同時にパケットを転送できる。図19に示すようにモジュール1(104)からモジュール2(105)の転送106とモジュール3(101)とモジュール0(102)との間の転送103は同時にできる。これが実現できれば、スループットを最大2倍に向上することができる。リングにおいてこれを実現した例は見あたらない。

【0044】従来の一般的な技術として、各モジュールから並列的に要求を集め選択する集中的アービトレーションの方式が知られる。この部分は、スイッチと同じトポロジーになる。集中的アービトレーション方式は共通部分が必要でリングの拡張性を減ずる。集中的アービトレーションをそのままリングで実現すると各モジュールからのパケット数は送信要求と同じになり、リングのスループットの大半を費やす。しかも送信要求を各モジュールが同時に行うには接続モジュール数倍の信号線が必要になる。

【0045】本リングの特性を活かして集中的アービトレーションをリングで実現することができる。ブロックパケットの転送には数サイクル掛かるので、送信要求をモジュール全体に対して1サイクルに1つ行う能力があればバランスが取れる。これにより送信要求を送るための送信権を獲得する機構を全体で1つ設ければ良いので、信号線はモジュール1つ分で良い。送信権を獲得したモジュールは、送信要求をアドレス0のモジュールに送る。アドレス0のモジュールは、各モジュールからの送信要求を同時に最大2つ選択して送信権を通知する。送信要求は、要求アドレス、宛先アドレスからなり、要求アドレスで送信権を通知する。信号を追加すれば、ブロックパケットの要求、ビジーの報告も行うことができる。基本リングの送達確認の信号線2本に加え並列転送のため、もう1組を追加する。

【0046】[並列送信要求選択方式] 従来の集中的アービトレーションでは1サイクルに1つの送信要求を選択している。本リングでは、1サイクルに2つの送信要求を選択すること、要求アドレス、宛先アドレスから転送可能性を判定することなど、従来技術にない処理能力が要求される。本発明のリングは、図20に示すように、送信要求をモジュール毎にバッファし1つ目の送信権を転送可能なブロックを並列的に判定してその中から最初のものを選択し、2つ目の送信権を最初の送信権と並行転送可能なブロックを並列的に判定してその中から最初のものを選択する高速論理方式である並列送信要求選択方式を用いている。

【0047】[ディレクトリ・リング] 各モジュール毎

に、それに属するキャッシュにあるラインの状態のみを格納するディレクトリを持ち、1カ所に集中してディレクトリを持たない分散方式がある。この分散方式では、各モジュールに属するディレクトリよりキャッシュの状態を集約し、自分のモジュールに属するディレクトリの状態を決め、メモリモジュールに通知する必要がある。このような場合の従来のキャッシュの状態制御は、状態の報告を各モジュールから受け、集約してメモリモジュールに配信する集中方式を用いている。この部分はトポロジ的にはスイッチと同じになる。集中的キャッシュ状態制御には共通部分が必要でリングの拡張性を減ずる。キャッシュの状態を調べる時間は一定でなく、キャッシュの状態の報告は、全モジュールから受け取ったことを確認する必要がある。従来の方式をそのままリングに適用すると全モジュールからキャッシュの状態をパケットで報告することになる。これでは、送信要求に対してそのモジュール数倍のパケットが必要になりリングのスループットの大半を消費する。

【0048】ディレクトリ・リングを用いることにより、本発明のリングの特性を活かしてキャッシュの状態制御をディレクトリ・リングで実現できる。報告を求めるモジュールから全モジュールについて報告を求めることができればパケット数は減らせる可能性がある。報告を求める間隔を報告の揃う時期に設定すれば高々1-2パケットで全モジュールからのキャッシュの状態の報告を受け取れる。パケットに各モジュール毎に状態を報告する領域を設けると、モジュール数倍の信号線を要する。キャッシュの状態の種類別に情報線を定めて集約できれば信号線が少なくても良いが、どのモジュールが報告したかの情報は失われる。

【0049】必要な機能は、全モジュールが報告したことを確認することである。この確認のためにキャッシュの状態を報告していないことを示す信号線を1本設ける。この信号線でキャッシュの状態報告できないことを知らせる。この信号線が有意でなければ全モジュールが報告したことになる。かくして、図21に示すようにキャッシュの状態報告を要求するための送信権の獲得、状態報告の要求、キャッシュの状態報告、未確定の報告、集約したキャッシュの状態の通知、確定の通知に必要な信号線を設けることによりリングで実現できる。未確定報告と同時に、その種類別に未確定原因を報告すれば、要求モジュールは効率よい再要求のタイミングを選択できる。

【0050】〔階層接続〕階層接続によりスループットを向上させる。ローカルメモリのキャッシュコヒーレンスを行うリングとして従来技術として前述したSCIがある。SCIにおいては、キャッシュコヒーレンスを全報知しないことを前提にパケットの問い合わせ、応答を階層的に行うためにパケット数が増え、確定までの時間が掛かる。また、ローカルメモリに対する完全なデ

ィレクトリを必要とするために容量も大きい。

【0051】本発明のリングは、ローカル・アクセスには性能の劣化がなく、リモート・アクセスには単純なディレクトリを設けて必要なモジュールのみにアクセスする手段を提供する。ディレクトリにエントリがないときにはリングの全報知の高性能を活かして全報知する方式を採用してレイテンシの改善と簡素化を行う。リングに属するローカルメモリのラインがリングの外で変更されているかが分かればリングの外にアクセスするかを決められる。リングに属するローカルメモリのラインがリングの外にコピーがあればコヒーレンスの動作をするかを決められる。リングに属するローカルメモリの全ラインに対して他のリングに変更有り、コピー有りの2ビットを持つディレクトリ（ローカル・ディレクトリ）を設ければ容量を余り必要とせずローカルアクセスについて他のリングをアクセスすることはなく性能劣化はない。変更或いはコピーのあるリングを特定する情報をディレクトリに追加すると必要な他のリングにのみアクセスすれば良いので性能は上がるがディレクトリの容量は大きくなる。ローカルメモリの全ラインに対してでなく最新のアクセスのラインを登録するディレクトリ（リモート・ディレクトリ）を設け、各リング或いは各リングのグループでの変更、コピーの状態のベクトルを持てば容量を小さくできる。変更のアクセスは排他制御のために頻発するのでディレクトリに優先して維持する。リモート・ディレクトリにエントリがあれば該当するリングだけにアクセスする。リモート・ディレクトリにエントリがなければ、全リングにアクセスする。図22に示す様にいずれの場合も、アクセスのパケットと兼用でき、1回のパケットの転送でコヒーレンスを取れるので時間は短い。キャッシュコヒーレンスの時間を短縮することでレイテンシ、全報知パケットを削減することでスループットをSMP並に改善することができる。ローカルメモリのアクセスの抑止は、ローカルメモリにもディレクトリを持つか、ローカルメモリを直接制御すれば、階層接続機構を通じて知らせる必要はない。その方法がないときには、ディレクトリが分散しているときの手段であるディレクトリリングを設ける必要がある。

【0052】〔ディレクトリのローカルメモリへの格納〕ディレクトリのローカルメモリへの格納により処理効率を向上させる。ローカル・ディレクトリの容量が大きくなるときには、ローカルメモリにローカル・ディレクトリを置きその写しを階層接続機構に持つ。階層接続機構が受信不能である時にディレクトリのアクセスのためにローカルメモリをアクセスするとその応答を受け取れない。このデッドロック状態を解決するために、従来技術である論理チャネル、或いはサブチャネルの概念を一部導入する。リングのアドレスを物理チャネル或いはチャネルと見立てて、リングの本来のアクセスとディレクトリのアクセスをそれぞれ論理チャネル0、1と解釈

する。この場合に必要になるのは、どの論理チャネルがパケットを受け取れるかであるので論理チャネルのステータスに相当するものとして受信不能の報告時に論理チャネル 0、1 を指定する応答信号を 1 本追加する。各モジュールは、フロー制御、ビジー状態をモジュールアドレスと論理チャネル 0、1 で管理する。論理チャネルの拡張、ステータスの追加には応答信号の本数を増やす。

【0053】[制御リング] 現状のリングでは、モジュールの障害検出、特定、切り離しはリングの通常動作を利用して行うので、そのための論理量が多い。障害モジュールの特定、切り離しのためには図 23 に示すように逆方向に回る 2 つのリングが必要になる。現状のリングは伝送系の障害が支配的である遠距離の伝送を前提とするため、この 2 つのリングは折り返し機能を実現する冗長系として効果がある。しかし、冗長系にするとレイテンシは、2 倍になりリングの構成要素は 2 倍になる。プロセッサバスではモジュールの障害に起因するリングの障害が支配的で、伝送系に起因する障害はむしろ例外的である。プロセッサバスでは、リングに接続するモジュールの障害に起因するリングの障害検出、特定、切り離し等のシステム運用が重要である。

【0054】本発明のリングは、障害検出、特定、切り離しの論理量を小さくして、構成変更、障害検出、特定、切り離しを独立して行うリング制御ユニットを設け、低速の制御データ線と制御クロックの 2 本の信号線とクロックからなる制御リングで各モジュールに接続する。リング制御ユニットからモジュールには転送によらずに配信する方式として、制御リングのモジュールの構成変更、障害に伴う各モジュールへの影響を軽減する。

【0055】図 24 に示すように、データ転送のための 1 つのリングと 3 本の信号線からなる制御リングにより、冗長リングを設けるよりも低コストで運用性の高いリングを構成できる。制御リングは直接的な制御をするので、各モジュールでの運用機能は不要になる。モジュールは障害を直接に報告できるので障害の検出の論理量が少ない。切り離しは障害モジュールにリングの中継を指示するだけなので切り離しのための論理量は少ない。制御リングにより、クロックの位相調整、受信線の遅延量をイニシャライズするので、各モジュールのイニシャライズ機能は簡素化できる。切り離しできない障害は、ケーブル、コネクタ、中継ゲート、ドライバ、レシーバに限定される。モジュールの挿抜に伴うリングの切断、復旧をモジュールから制御リングを経由して行うので時間を短縮できる。モジュールの内部には、イニシャライズ中もクロックを供給してイニシャライズ終了後位相調整されたクロックに連続的に切り替えるので各モジュールは連続動作できる。

【0056】[構成変更の簡便性] 従来のリングにおいて、中継部分を含む接続装置の追加、除去はケーブルの接続変更を要する。プロセッサリングにおいては、装置

の追加、除去に相当するモジュールの挿抜を行うには、短絡スイッチとの入れ替えになる。その際、モジュールの引き抜きと短絡スイッチの挿入と 2 つの動作が必要である。本発明では、冗長リングを用いずに構成変更の簡便性を上げる。モジュールの挿抜に伴い、その力を使用して連動して挿抜する短絡スイッチを考案してプロセッサバスへの適用に伴う構成変更の時間の短縮と簡便性を実現した。

【0057】モジュールの挿抜を動作中に前触れなく行うとパケットの送信と応答など一連の動作が中止されシステムは異常に停止する。一般には、モジュールの挿抜を行う前にオペレータによりシステムを停止し、挿抜後にシステムを立ち上げるのでシステムの停止時間が長い。そこで、モジュールの挿入を事前に検知するセンサを設け、次のモジュールにセンサを接続して制御リング経由で通知する機能を設ける。また、モジュールの引き抜きを事前に検知するセンサを設置し制御リング経由で通知する機能を設ける事前に挿抜するモジュールだけを切り離して、その後は必要最小限の時間だけ自動的にシステムを停止すれば良いので、業務処理の停止時間を短縮することができる。

【0058】[論理的順序づけ、双方向リング] 本発明のリングは、メモリを用いるセマフォアのための論理的順序づけを保証できる。また、双方向リングを使えば、スループットを約 2 倍に向上できる。

【0059】[レイテンシ] 図 25 に本発明のリングによるレイテンシの改善内容を示す。

【0060】[動作周波数] 図 26 に本発明のリングによる動作周波数の改善を示す。

【0061】[増設性] 本発明のリングは、同一周波数、リング配線での送信権リングの追加、階層接続により性能を向上することができる。また、同一周波数、リング配線での双方向リング化による性能向上とクラスタ化が可能である。クラスタ化にいたる増設性のパスと同一の LSI 技術を用いた性能向上の一例を図 27 に示す。

【0062】[単一リングの性能] 動作周波数の限界は、送信ラッチの出力毎にスキュー補正をしない場合は、スキュー補正回路の精度によるスキューと送信ラッチのスキューの重量で決まる。送信ラッチの出力毎にスキュー補正をすればスキュー補正回路の精度だけで決まる。ドライバの周波数特性は一般にこれより大きいので、概ねドライバの動作周波数から考えて 400MHz 程度と考えられる。

【0063】レイテンシを 8 モジュール、配線長の最大 10cm で考える。中継モードのモジュールのレイテンシは、ゲートディレイ 2 段 (0.3ns)、ドライバの遅れ (0.8ns)、配線 (10cm) を含む実装遅れ (1ns) とスキューの平均遅れ (0.4ns) とで、合計 2.5ns である。転送モードのモジュールの



レイテンシは、中継モードのモジュールのレイテンシに加え、ラッチの受信ウィンドの中央（平均0.4サイクル、400MHzの時、1ns）とラッチの遅れ（0.5ns）で、合計4nsである。

【0064】中継モードのモジュールが6、転送モードと回送モードが1であると、リング1周の平均レイテンシは23.5nsになる。リング一周のサイクル数は、400MHzの場合10サイクルになる。250MHzの場合、転送モードと回送モードのレイテンシがそれぞれ0.8ns増えるのでリング1周の平均レイテンシは25.1ns、リング1周のサイクル数は7サイクルになる。モジュール間の平均レイテンシは、ほぼリングの半周になるので400MHzの場合、約12.5ns、250MHzの場合14nsになる。

【0065】[階層接続の性能] 階層接続では、（1）送信副リングでの送信モジュールから接続機構、（2）送信接続機構の処理、（3）主リングでの送信接続機構から受信接続機構、（4）受信接続機構の処理、（5）受信副リングでの接続機構から受信モジュールまでがレイテンシになる。

【0066】（1）、（3）、（5）は、単一リングのレイテンシと同じであるので単一リングのレイテンシの3倍になり、250MHzの場合42nsになる。

（2）の送信接続機構の処理は、受信と送信処理である。ディレクトリのメモリへのアクセスはない。受信と送信にそれぞれ6サイクル、250MHzの場合24nsになる。全報知のパケットは、ディレクトリをアクセスするのでこれを加える。（4）の受信接続機構の処理は、ディレクトリのアクセスと副リングへのパケットの転送を並行してできるので受信と送信に6サイクル、250MHzの場合24nsとなる。これらの結果、全体では90nsになる。

【0067】次に、本発明の実施例について説明する。

【0068】[基本リング] 最初に、基本リングについて説明する。

【0069】[接続形態] 対象とするプロセッサ間リングインターフェースは複数のプロセッサモジュール（PM）、入出力モジュール（IO）、メモリモジュール（MM）、リング制御ユニット（RC）等を接続する。PM、IO、MM合わせて8モジュールをリングに接続する。PM、IO、MMの各モジュールの信号は、基本リングとして選択線（1本）、アドレス線（3本）、データ線（72本）、制御リングとしてクロック（1本）、制御クロック（1本）、制御データ線（1本）である。各モジュールは、基本リングの選択線、アドレス線、データ線を前のモジュールから受信し、次のモジュールに転送する。リング制御ユニットは、基本リングを動作させる高速のクロック、基本リングを制御する低速の制御クロック、制御クロックに同期した制御データを発生し配信する。

【0070】モジュールには、回送モードと転送モードの動作モードがある。リングに1つのモジュールは選択線、アドレス線、データ線をクロックに同期して転送する回送モードで動作する。基本リングの1周は、クロックの整数倍の遅延で2サイクル以上に設定する。図28に基本リングの各モジュール間の接続を示す。

【0071】各モジュールは、制御リングの制御クロックで制御データを処理し、自分が宛先であるときに受信する。内容に応じて、モジュールの状態についての応答を制御データに乗せて送信する。リング制御ユニットは、各モジュールにポーリングして各モジュールの情報を受信する。クロックは、動作周波数の2分の1の周波数を配信してモジュールで動作周波数に倍周する。クロックの動作周波数の制約は軽減される。以下の説明では簡単のため動作周波数のクロックを供給することで説明する。図29にリング制御ユニットに関するリングインターフェースの信号を示す。

【0072】[送信権の先行獲得] パケットは要求元、要求番号及び種別等のヘッダと本体とからなるが、その形式及び内容についてはリング固有の課題はないのでここには述べない。送信権については、選択線、アドレス線を受信しその内容に従い獲得する方式がよく知られる。図30に従来の方式を適用する例を示す。選択線720は入力ラッチ721に格納され、送信要求722について送信権処理724を行い、結果を出力ラッチ723に格納する。レイテンシは最低1サイクルは生じる。このレイテンシを短縮する方式をここに述べる。

【0073】図31に、選択線をゲート2段の遅れで転送できる方式を示す。送信要求727はこのまま選択線720と論理和728をとり転送する。選択線720の内容と送信要求727をしていることから送信権獲得確認726をする。選択線は、自分の送信権を放棄（725）する時には選択線を論理積729で抑止する。

【0074】図32に受信線についての従来の方式を適用する例を示す。受信線733は受信ラッチ730に格納され、送受信処理731を行い送信ラッチ732に格納する。レイテンシは、最低1サイクルは生じる。

【0075】図33に、受信処理734と送信処理735を並行して行う方式を示す。受信線733は受信処理（734）され、並行して送信処理735を行い、送信の場合は論理積737により受信線を抑止する。

【0076】パケットの宛先はアドレス線を用い、パケットはデータ線を用いる。なお以下の説明で明確になるが、アドレス線は説明のため分離しているが、実際にはパケットのヘッダに含まれるので実際の信号線はデータ線に含まれる。データ線（72本）の内8本はパケットの内容に応じてエラー訂正符号或いはパリティに用いられる。実際のデータは、64本（64ビット）である。宛先とパケットはリング上のクロックに同期した1サイクルで転送する。選択線は2サイクル後のデータ線、ア

ドレス線にパケットと宛先があることを示す。

【0077】図34は、時間2(201)で選択線を1(202)にしてその2サイクル後の時間4(203)でアドレス線に宛先3(204)を、データ線にパケットA(205)を、時間4(203)で選択線を1(206)にしてアドレス線に宛先5(207)を、データ線にパケットB(208)を送信することを示す。

【0078】[受信]信号の変化の後縁でラッチからラッチに転送する受信方式を図35を用いて説明する。受信10は当該モジュールが送信をせず受信状態にあることを示す。前のモジュールからのデータは、データ線11を経由して当該モジュールに転送される。受信10が1であるので論理積12が成立し、データ線11は論理和回路13を経由して次のモジュールに転送する。データ線の内容は受信ラッチ14に常にラッチする。アドレス線についても同様に受信ラッチ15にラッチし当該モジュールの自分アドレス16と一致回路17で比較する。

【0079】選択線の内容が1であれば、2サイクル後のアドレス線とデータ線に宛先とパケットがあることを示す。選択線はラッチ18、ラッチ19、ラッチ20に1サイクルづつ遅れてラッチされる。ラッチ20の選択線の内容は受信ラッチ15より2サイクル進んでいるので選択線の内容1でアドレスの一致論理17が1ならば論理積21が成立して、論理積22を経由して受信パケット23として受信する。図36に受信動作をタイムチャートに示す。

【0080】[送信権獲得と送信]送信したいパケットがあれば、選択線を1にする。受信した選択線が0であれば送信権を獲得し2サイクル後に送信モードとなる。受信した選択線が1であれば、送信権が獲得できなかったことを知る。そのまま受信を続け、かつ送信権の獲得待ちになり選択線が0になるのを待つ。選択線が0になれば、2サイクル後に宛先とパケットをアドレス線、データ線に乗せる。送信したパケットがリングを1周すれば選択線を0にし送信権を放棄する。そのパケットがリングを一周する所定のサイクルより2サイクル前迄に送信要求があれば引き続いて選択線を1にして送信権を獲得できる。

【0081】図37に送信動作の概念図を示す。送信要求600があれば1(601)を出力する。受信した選択線602が1(603)ならば引き続き送信要求600を1(601)にする。受信した選択線602が0(604)ならば送信要求を0(605)にして送信権を獲得したことを知る。選択線の出力(606)は、受信した選択線602と送信要求600の論理和を出力する。

【0082】図38を用いて送信動作を述べる。当該モジュールが送信したいことは送信30を1にすることで示されている。以前に送信30を要求してそれが処理さ

れていれば、要求(ラッチ)31は0である。論理積32が成立して受付(ラッチ)33が1となる。受付33は、送信30が受理されたことを送信受付(34)信号でモジュールの送信要求論理に知らせる。ついで新たな送信要求を調べ送信30信号に伝える。毎サイクル送信要求をするには、受付33の結果を用いて1サイクルで更新する必要がある。更新回路は、リングインターフェースの問題ではないのでここでは述べない。論理積32が成立することにより同じく要求31(ラッチ)が1になる。要求31(ラッチ)の内容は直ちに論理和35を通じて次のモジュールに伝える。当該モジュールが送信権を獲得後リング1周したサイクルであれば、送信権放棄36は1になり論理積37が成立せず選択線の転送を抑止する。

【0083】他のモジュールが送信権を獲得していないサイクルならば選択線の内容は0であるので単に要求31の内容を伝える。他のモジュールが送信権を獲得したサイクルならば選択線は1でありかつ論理積37が成立してその送信権と当該モジュールの送信要求は論理和35に重畳される。受信した選択線が0であれば、他のモジュールが送信権を獲得していないことを示し論理積38が成立し獲得39(ラッチ)が1になり送信権を獲得したことを知る。送信権放棄36が1であるときは送信権の処理は選択線が1であっても選択線が0と同じ動作をする。送信権放棄36が1であれば、当該モジュールの送信パケットがリングを1周したことを合い召し論理積42が成立し獲得39(ラッチ)が1になり送信権を獲得したことを知る。選択線が0になってから2サイクル後、獲得39(ラッチ)が1になってから次のサイクルに送信モードとなりデータ線、アドレス線にパケットと宛先アドレスを転送する。受信した選択線が1であれば、論理積38が成立せず獲得39(ラッチ)は0であり送信権を獲得できなかったことを知る。送信権獲得待ちで選択線が0になると論理積38が成立して獲得39(ラッチ)が1になる。選択線が1で、要求31(ラッチ)が1で送信権放棄36が0のときには論理積40が成立して要求を出し続ける。送信権の獲得待ちになり選択線が0になるのを待つ。選択線が0、要求31(ラッチ)が1、且つ送信30が1ならば論理積41が成立して引き続いて送信権を獲得することになる。送信権放棄36が1、要求31(ラッチ)が1、且つ送信30が1ならば論理積43が成立して引き続いて送信権を獲得することになる。論理積41、論理積43が成立すると受付33(ラッチ)が1になり送信30を受付けたことを知らせる。獲得39が1になると論理積47が成立して送信ラッチ48に宛先アドレスをラッチし論理和50を経由して次のモジュールに送信する。獲得39が1になると論理積51が成立して送信ラッチ52に送信パケットをラッチし論理和53を経由して次のモジュールに送信する。ラッチ44が1になるとアドレス線の内容を次

のモジュールに転送しないように論理積54で抑止する。ラッチ44が1になるとデータ線の内容を次のモジュールに転送しないように論理積55で抑止する。送信論理は受信論理と同じく2段のゲート論理で構成される。

【0084】図39に送信動作の一例を示す。送信30が1になり、要求31が0であることから要求31、受付33が1になる。受付33が1になることにより、送信30が更新され1サイクルの中で0になる。選択線

(入力)が引き続き1であるので要求31は1である。受付33は送信30が0であるので次のサイクルは0になる。選択線(入力)が0になる次のサイクルで獲得39は1になる。獲得39が1になる次のサイクルでラッチ44が1になる。ラッチ44により送信ラッチ48が更新しアドレス線に出力する。ラッチ44により送信ラッチ52が更新しデータ線に出力する。リング1周のサイクルである6サイクル後に送信権放棄36が0となる。当該モジュールで送信要求が引き続かないので次のサイクルで選択線(出力)は0となる。

【0085】[スキュー補正]従来方式においては、クロックを受信線に対してスキューの設計値にセットアップ時間を加えただけ遅延させ、パルス幅はクロックのスキューがスキューの設計値だけ遅れても受信できる必要がある。動作周波数は、スキューの設計値の2倍にセットアップ時間を加えて決められる。図40にこれを示す。

【0086】図41に受信線に合わせてクロックの位相を設定する方式を示す。動作周波数はスキューの設計値にセット時間を加えたものになる。レイテンシは、スキューの実績値とセット時間の和に短縮される。図42に受信線毎にスキュー補正を行う方式を示す。動作周波数は、スキューの補正精度で決まる。

【0087】[受信位相での処理]従来方式ではリングのクロックと各モジュールのクロックの位相は独立である。平均すれば2分の1のサイクルの遅れを生じる。各モジュールは、受信線の位相に合わせて動作すればレイテンシを削減できる。リングに1つのモジュールはリング1周するサイクルをサイクルの整数倍にする必要がある。図43に従来方式を示す。リングのクロック740によって受信線733を受信ラッチ730に格納してモジュールのクロック742によりサイクル調整を行いモジュールのサイクルに同期する。送受信処理731を行い送信ラッチ732に格納して転送する。

【0088】リングのクロックに同期して動作する方式を図44に示す。受信線733は、クロック742を位相調整(740)して受信ラッチ730に格納し送受信処理731を行う。受信線733と送信内容を論理積737、論理和736で選択しリングクロック740で転送ラッチ743に格納する。各モジュールは、転送ラッチ743の出力を次のモジュールに転送する。これを転

送モードと呼ぶ。リングに1つのモジュールは、転送ラッチ743をサイクル調整741を経てクロック742で回送ラッチ744に格納して次のモジュールに転送する。この同期方法を回送モードと呼ぶ。

【0089】[転送モードの動作]選択線、アドレス線、データ線等の受信線のラッチの位相を受信線の遅延とスキューに合わせて転送ラッチのクロックの位相を定めて動作周波数、レイテンシを改善する。動作周波数には、クロックのスキューの設計値と受信線のスキューの設計値の和が直接的に影響を与える。クロックのスキュー補正を行えば動作周波数に与える影響はクロックスキューの設計値を加えずに、受信線のスキューの設計値だけでよいので2分の1にすることができる。クロックのスキュー自動補正を行えば、リングのレイテンシに与えるスキューの影響は、クロックのスキューの設計値を加えずに、受信線のスキューの設計値でなく受信線のスキューの実現値だけでよい。一般にスキューの設計値と実現値との乖離は大きいので動作周波数の改善に寄与する。

【0090】図45に転送モードの送信動作の概念図を示す。最も早く受信した受信線(610)と最も遅く受信した受信線(611)の中間の2分の1サイクル後に遅延クロック(612)を設定する。受信したデータは遅延クロック(612)で転送ラッチに格納して次のモジュールに転送する。受信した選択線(613)の位相は一般に両者の中間にある。受信した選択線613はラッチ(614)する。

【0091】送信要求615があれば1(616)にする。ラッチ(614)した選択線613が1(617)であるので引き続き1(616)にする。ラッチ(614)した選択線613が0(618)になると送信要求615を0(619)にする。選択線620にはラッチ614と送信要求615の論理和を出力する。

【0092】[位相調整回路]遅延クロックは、図12に示すように倍周されたクロックから位相調整回路を経由して作る。図46はゲートの遅延を利用した位相調整回路であって、半導体の特性ばらつきを考え、最小の遅れでも2分の1サイクルの遅延を保証できるように遅延量を取る。遅延量は遅延設定レジスタで指定する。遅延量は、選択回路1段、選択回路2段、選択回路2段にゲート2段からゲート16段を加える8段階で合計10段階、負クロックの入力をいれると20段階になる。但し、半導体の特性ばらつきにより最大の遅れでは、クロックと負クロックの遅延が重複することがある。例えばゲート8段階で遅延が2分の1サイクルに達すると6段階になる。

【0093】[調整クロック]調整クロックは、4サイクルの間位相を連続して変化する信号である。前のモジュールより設定用のパターンを送信して遅延設定を行う。前のモジュールは、信号の変化点を容易に特定でき

るように選択線、アドレス線、データ線に1000の4サイクルの繰り返しパターンを送出する。設定用パターンの4サイクルの一つある信号の変化点を特定するために4サイクルの間連続的に位相を変化する調整クロックを設ける。調整クロックは、4サイクルの周期で1サイクル単位に位相を変化させる機構、位相調整回路、遅延がそのサイクルにあることの範囲検出回路を接続して発生する。1サイクル単位に位相を変化させる機構は、リングカウンタで4サイクルの周期の1000の繰り返しパターンを発生し、サイクルを選択する機能を有する。上記信号とクロック及び逆極性のクロックとの論理積を取りどちらかを選択することで2分の1サイクル単位に位相を変化させる。上記の信号を位相調整回路の調整クロックに入力する。範囲検出回路は、位相調整回路の出力である遅延クロックと遅延前の調整クロックとの論理積を取りセットリセットラッチに入力してその出力信号を判定する。出力信号が1ならばその遅延の範囲が2分の1サイクル以下であると判定する。遅延がその2分の1サイクルの範囲を超えるか、位相調整回路の遅延量が最大になれば2分の1サイクル先の位相に移る。かくして、4サイクルの周期の可変遅延を有する受信位相決定用の調整クロックを作成できる。

【0094】[受信位相の決定] 受信位相の決定の手順を以下に示す。

【0095】すべての受信線について設定用パターンの0を受信する調整クロックの位相を定める。少なくとも1つの受信線について1を受信する遅延クロックの位相を定め、前縁とする。すべての受信線について1を受信する遅延クロックの位相を定め後縁とする。前縁と後縁の中間の2分の1サイクル後に受信位相を定める。前縁と後縁の中間は次の手順による。前縁の検出時に計数を開始して後縁で終了するスキューカウンタを設ける。次に前縁に調整クロックを設定する。スキューカウンタの2分の1迄調整クロックの位相を進めて前縁と後縁の中間に位相調整回路の遅延量を定める。前縁と後縁の中間の2分の1サイクル後のクロックは、調整クロックを作成する元のクロックとは逆極性のクロックを位相調整回路に入力することで得られる。

【0096】[中継モード] 中継モードでは、転送ラッチの入力を受信データとして転送する。転送ラッチの入力の切り替わりは転送ラッチのタイミングの2分の1サイクル前である。転送モードの送信では、送信内容は転送ラッチのタイミングで切り替わっていて、1サイクル後に転送ラッチにラッチしてから転送している。中継モードの送信では、1サイクル前の送信内容を転送ラッチの入力の切り替わりに合わせるため2分の1サイクル遅延させて中継する必要がある。このため、送信内容を転送ラッチの逆極性のクロックでラッチして2分の1サイクル遅延させて転送する。

【0097】[回送モード] 回送モードでは、転送ラッ

チの内容を回送ラッチに移せばよい。但し、転送ラッチのタイミングの位相が回送ラッチのセットアップ時間内にあれば、いったん調整ラッチに移してから回送ラッチにラッチする。モジュールの接続数が少ない場合にリング1周のサイクル数を2サイクルにするため、或いは後述するように2つのリングのリング1周のサイクル数を同一にするために調整ラッチ1を設ける。回送ラッチのラッチタイミングの前後の内容が同じであることで、転送ラッチと回送ラッチのラッチタイミングの重なりのないことを確かめる。回送モードのモジュールは、10の繰り返しパターンを送信する。回送ラッチのタイミングを規定値だけ遅延させて転送ラッチの出力をラッチする。これは回送ラッチのタイミングの後の転送ラッチの出力を見ることになる。転送ラッチの出力を転送ラッチのタイミングを規定値だけ遅延させて、回送ラッチのタイミングでラッチする。これは回送ラッチのタイミングの前の転送ラッチの出力を見ることになる。両者の信号が一致すれば、回送ラッチと転送ラッチのタイミングの重なりがないので調整ラッチを経由する必要はない。

20 【0098】図47を参照してサイクル調整回路を説明する。転送ラッチと回送ラッチのタイミングに重なりがあると、調整ラッチ使用311が1になる。モジュール制御ユニットは、リング制御ユニットの指示により調整ラッチ1使用312を1にする。転送ラッチ313は調整ラッチ使用311が1ならば調整ラッチ314に移す。調整ラッチ1使用312が1ならば調整ラッチ314の内容を調整ラッチ1(315)に移す。転送ラッチ313の内容を直接、調整ラッチ314を経由、調整ラッチ1(315)を経由するかを論理積316、論理積317、論理積318により選択して回送ラッチ319に移す。

30 【0099】[受信線毎のスキュー補正] 受信線のスキューの設計値は、動作周波数に影響を与える。選択線、アドレス線、データ線からなる受信線毎に可変遅延回路を設けスキューの自動補正を行えば動作周波数に与えるスキューの影響は補正精度だけになる。動作周波数は、ドライバ、LSI・基板内の配線、パッケージ、コネクタ等の伝送特性で定まる。可変遅延回路は、ゲートの遅延を利用した位相調整回路と同様のものでよい。但し、反転入力、調整用の入力信号は不要である。最小の遅延量はリングの設計値で決まり2分の1サイクルまでの必要はない。

40 【0100】図48に受信線毎にスキュー補正する転送モードの概念図を示す。受信位相の時と同様に位相設定用の1000のパターンを前のモジュールより送信する。受信線毎の可変遅延回路に遅延量を0にして最も遅い位相の受信線611に調整クロックの位相を合わせる。受信線毎の可変遅延回路の遅延量を一齐に増やし、受信線毎に転送ラッチで1を検出する最後の可変遅延回路の遅延量を格納する。遅延の走査完了後一齐に格納し

た遅延量に戻す。この結果、すべての受信線の転送ラッチの入力は、同じ位相になる。調整クロックの元のクロックの逆極性のクロックを位相調整回路に供給することにより2分の1サイクル遅延させ受信位相のクロックとする。

【0101】〔全報知の packets〕 packets を各モジュールに同時に送信するために、全報知信号（1本）を設ける。全報知信号はアドレスと同時に1にすることにより全報知を示す。全報知信号が1の時にはアドレス信号に要求元のモジュールのアドレスを示す。全報知信号は、説明を簡潔にするために信号線として分離しているが以下の説明で明らかになるように実際には packets のヘッダに含められるので信号線の本数を増加させない。

【0102】〔フロー制御〕 リングのレベルで packets の送達確認を packets によらず2本の信号線の追加で高速に行い、複雑な制御回路を回避し、スループットの低下を除去する手段を設ける。全報知の packets に対しての送達確認の手段も設ける。リングのレベルでの再送信を減らすために、トラフィックに応じたフロー制御のアルゴリズムを与える。

【0103】〔送達確認〕 図49により送達確認の説明をする。受信モジュール400が、受信の2サイクル後に受信不能402を送信モジュール403に報告する信号（1本）401を設ける。送信モジュール403が受信不能の報告を受けたことを報告受領後2サイクル後に各モジュールに通知する信号（1本）404を設ける。受信不能を報告した受信モジュール400は、受信可能になったことを全報知の packets で送信する。送信モジュール403は、2サイクル後の送信 packets （408）に加え、その2サイクル後の受信不能信号401、受信不能信号401を送信モジュール403が受信してから2サイクル後の受信不能通知信号404の送信権を獲得する。

【0104】〔受信不能〕 送信モジュール403は packets 送信の2サイクル後に受信不能信号401を0（405）にする。受信モジュール400は、 packets を受信できなかった時に受信不能信号401に1（402）を入れる。 packets を受信できた時には何もしない。送信モジュール403は、受信不能信号401が0（405）のままであることで送達確認をする。受信できれば何も受信不能信号401にしないので0のままになり全報知に対しても送達確認できる。

【0105】〔受信不能通知〕 送信モジュール403は各モジュールに受信不能404を通知する。受信不能信号401が1（402）ならば、受信の2サイクル後に受信不能通知信号404を1にする。受信不能信号401が0（405）のままならば受信の2サイクル後に受信不能通知信号404を0にする。各モジュールは、特定の受信モジュールに宛てた packets が受信不能である場合は、特定の受信モジュールが受信不能であることを

知る。各モジュールは、全報知の packets が受信不能である場合は全報知 packets が送信不能であることを知る。送信モジュールに到達する前の受信不能信号401が1（402）であることをスヌープ（406）できる受信モジュール407はそれで受信不能を知ることができるので処理に利用しても良い。

【0106】〔受信不能原因の報告〕 受信不能線に加え、原因別の受信不能原因線を追加して受信不能原因を報告してより強力なフロー制御を行うことができる。同様に受信不能通知線に加え、原因別の受信不能原因通知線を追加して、より強力なフロー制御を行うことができる。

【0107】〔再送信〕 特定の受信モジュールに宛てた packets が受信不能になると各モジュールは、そのモジュールへの送信及び全報知の送信を抑止する。全報知の packets が受信不能になると各モジュールは、全報知の送信を抑止する。受信不能を報告した受信モジュールの受信バッファの空きエントリがリングのサイクル数の2倍以上になると受信可能を全報知の packets で送信する。各モジュールは、その受信モジュールが受信可能になったことを知る。受信不能の受信モジュールに宛てた packets の送信を再開する。送信不能である全報知の packets の再送信は、複数の受信モジュールが受信不能であることがあり得ることから特定の受信モジュールに宛てた packets の後にする。全報知の packets の再送信は、空きがリング1周のサイクル数連続するか規定の時間送信を抑止する。全報知の packets の受信不能が増加すれば、送信を再開する迄の規定の時間を長くする。全報知の packets の受信不能がなければ、送信を再開する迄の規定の時間を短くする。

【0108】〔フロー制御〕 フロー制御は、自分の packets の送信権の放棄により行う。次に接続されているスロットは最高の順位の送信権を得る。空き packets が、連続して規定値以上のサイクルの間に存在しない場合は引き続き送信する packets があっても送信権を放棄する。規定値は、平均のビジー率が低ければ自動的に大きく設定する。バースト的なビジー状態でのビジー率の低下を防ぐ。平均のビジー率が高ければ、特定の受信モジュールに宛てた packets を優先し、ビジー率が低ければ全報知の packets を優先する。

【0109】〔複数サイクルの packets〕 packets を複数サイクルで転送することにより、選択線を中心とする各モジュールの論理動作のサイクル数が増加し、またモジュール間の信号線の本数の削減によりコストの低減を図れる。2サイクルの packets に対しては2分の1のクロックを供給しているので配信される2分の1のクロックの位相と packets の最初のサイクルの関係を各モジュール毎に固定する。2サイクル以上の場合は、 packets のサイクル数の周期の packets クロックを各モジュールで作成する。 packets クロックは、最初のサイクルが1

で残りのサイクルが0である。各モジュールでのパケットクロックの同期は、リング制御ユニットから同期指示を出し、回送モードのモジュールはパケットのサイクルの最初を1残りを0にする繰り返しパターンを同期パターンとして送信する。各モジュールは、特定の同期パターンを受信してパケットクロックを同期させる。送信権の制御は、選択線を使用してパケットクロックが1であるサイクルに行う。3サイクル目から選択線は、0にされる。これにより、パケットクロックの同期状態を各モジュールで監視する。後述するディレクトリ・リング或いは選択リングは、複数サイクルに分けて送信するので2サイクルのパケットであれば例えばそれぞれ11本が6本に、24本が13本になる。

【0110】[ブロックパケットの転送] キャッシュラインの様な固定長のメッセージを送信する場合、連続的にパケットを送信することが保証されていれば、パケットへの分割転送の無駄が省けてスループットが向上し、受信モジュールでのメッセージの組立が容易になる。ブロックパケットの送信を行う時には、あらかじめブロックパケットの送信を通知して通常のパケットの抑止を行う。ブロックパケットの送信権の獲得は、ブロックの最初のサイクルで行う。ブロックパケットの送信権は、リング1周のサイクル数がブロックパケットのサイクル以下ならばブロックパケットのサイクル数の後に引き継ぐ。ブロックパケットの送信権は、リングの1周のサイクル数がブロックパケットのサイクル数以上ならばリング1周のサイクル数後に引き継ぐ。通常のパケットには、リング1周のサイクル数だけの独立した送信権がある。リングの1周のサイクル数がブロックパケットのサイクル数以上ならばその端数を除いた倍数だけの独立した送信権がある。ブロックパケットの送信をすることを予約し、それを各モジュールに通知して通常のパケットの送信を抑止してからブロックパケットの送信を行う。

【0111】図50を参照して独立した送信権が2つの場合を説明する。ブロックパケットを送信する予約権を獲得するために、予約信号(1本)を設ける。ブロックパケットの予約権はサイクル毎に存在しないので、その存在を最初に明示的に示すために起動する。これは1つの信号からなるトークンを発行することを意味する。パワーオン時にアドレス0のモジュールは、予約信号を1サイクルの間1(予約1(450))にして予約権を1つ起動して次のサイクルから単に中継する。最初の予約信号の送信後、リング1周のサイクル数にいたる迄にブロックパケットのサイクル数があれば独立した送信権を必要数(予約2(451))起動する。ブロックパケットを送信したいモジュール(452)は、予約信号に0を送信し続ける。ブロックパケットを送信したいモジュール(452)は1(予約1(450))の予約信号を受信すると予約権を獲得し、2サイクル後からは予約信号には0を送信せず単に中継する。ブロックパケットの

サイクル数は2以上なので次の予約信号(予約2(451))には影響がない。ブロックパケットの予約権を通知するため予約通知線(1本)を設ける。予約権を獲得の2サイクル後から予約通知線に1サイクルの間1(予約通知1(453))を乗せ、リング1周後のブロックパケットのサイクル数分(この場合4)の送信権の予約を各モジュールに通知する。予約権を獲得したモジュール(452)はリング1周のサイクル後予約通知線に1サイクルの間0(予約3(454))を送信する。ブロックパケットの予約権を獲得したモジュール(452)は、予約通知の2サイクル後に選択線を1にして送信要求(送信要求1(455))を行い送信権を獲得する。ブロックパケットのサイクル数の間も同様に送信権を獲得できる。

【0112】次の予約権はブロックパケットのサイクル数がリング1周のサイクル数より大きいならば予約した後のブロックパケットのサイクル数の後になる。この場合はブロックパケットのサイクル数がリング1周のサイクル数以下なので、次の予約権はリング1周のサイクル後になる。ブロックパケットを引き続いて送信しないならば、予約線を1サイクルの間1(予約3(454))にして予約権を放棄する。次のサイクルからは1を送信せず単に中継する。引き続いてブロックパケットを送信したい場合は、何もせずに予約権を獲得する。各モジュールは、予約通知の受信の2サイクル後からブロックパケットのサイクル数の間通常パケットの送信要求を抑止する。

【0113】ブロックパケットの送達確認は、ブロックパケット受信不能信号、ブロックパケット受信不能通知信号を設けて行う。ブロックパケットが受信不能ならば、最初のパケットの2サイクル後にブロックパケット受信不能信号を1にして報告する。送信モジュールは、受信不能信号2サイクル後にブロックパケット受信不能通知信号を1にしてブロックパケットの送信不能を通知する。送信モジュールはその受信モジュールに対するブロックパケットの送信を抑止する。受信不能を報告したモジュールは、全報知パケットによりブロックパケットが受信可能になったことを通知する。通常のパケットの送達確認と同様に、原因別の受信不能報告及び通知の機能を拡張できる。通常のパケットとしての送達確認は、受信不能信号、受信不能通知信号を使用する。平均ビジー率が定められた値より大きければ連続して送信権を獲得しない。これにより送信権を獲得した次のアドレスのモジュールが優先的に送信権を獲得できる。予約権を保有しているモジュールは明示的でないのでアドレス0のモジュールで監視を行う。アドレス0のモジュールはフロー制御と予約信号の正常性の監視のため、予約信号がブロックパケットのサイクル数或いはリング1周のサイクル数のどちらかの周期以上で受信すること且つあらかじめ定められた一定時間以内に受信することを確認す

る。定められた数のブロックパケットを連続して送信することは、予約線、予約通知線と同じ働きをする長ブロック予約線、長ブロック予約通知線を追加することで実現できる。長ブロックパケットを複数個纏めた超パケットを送信するには、長パケット予約線に相当する超パケット予約線、長パケット予約通知線に相当する超パケット予約通知線を設ける。同様に、送信権獲得のための信号を追加することにより階層的な構造の異なる長さを持つ複数種類のブロックパケットを送信できる。図51にリング1周のサイクルより大きなブロックパケットの転送の例を示す。

【0114】[ディレクトリ・リング] ディレクトリ・リングにより、モジュールのキャッシュにデータを保有していることを送信モジュールに報告し、送信モジュールがそれを各モジュールに通知することができる。送信モジュールは送信権を獲得して報告を要求し、報告に基づいてそれを通知する。図52にキャッシュにある報告がある場合、図53にキャッシュにない場合、図54にその報告の要求では確定せず、引き続き報告要求をする場合を示す。

【0115】ディレクトリ・リングの送信権の制御を行う要求線500、読み出しパケットに対応する要求番号を乗せる要求番号信号501(6本)、キャッシュラインの存在の有無を確認中であることを報告する未確定信号502(1本)、キャッシュラインがあることを報告する保有信号503(1本)、その結果を通知する確定信号504(1本)、保有通知信号505(1本)を設ける。信号線を要約すれば、11本の追加になり、2サイクルでパケットを転送する場合には6本の追加になる。先に述べるディレクトリのローカルメモリのアクセス待ちの様な未確定の状態が完了する迄に多くの時間を要することを通知する長期未確定信号(1本)を設けることもできる。再要求を行う時間を待たせて、無駄な要求を減らすことができる。要求線500により2サイクル後の要求番号信号501、要求線500より4サイクル後、要求番号信号501から2サイクル後の未確定信号502及び保有信号503、要求線500よりリング1周のサイクルに6サイクルを加えたサイクル後の確定信号504及び保有通知信号505の送信権の制御を行う。報告を要求するモジュールは、送信権を獲得(506)して2サイクル後にその要求番号を要求番号信号(507)に乗せる。引き続いて、未確定信号502、保有信号503を0(508)にする。キャッシュ状態を報告するためモジュールは、要求線が1(506)になっている2サイクル後の要求番号(507)を読む。モジュールは、その要求番号(507)の示すラインがキャッシュになれば、何もしないでリング1周のサイクルの2サイクル後の確定信号(504)を待つ。モジュールは、その要求番号(507)の示すラインがキャッシュにあれば、保有信号(503)を1(509)に

し処理を完了する。

【0116】図52の未確定信号(502)のx(510)は、0あるいは1を示す。図52、図53、図54において、リング3周の0サイクル目の受(511)は送信モジュールが受信することを示す。モジュールは、その要求番号(501)の示すキャッシュラインの存在の有無の確認中は、未確定信号(502)に1(図54の512)を入れ、リング1周のサイクルの2サイクル後の確定信号(図54の513)を待つ。要求モジュールは、リングの1周のサイクル後の未確定信号(502)と保有信号(503)を受信する。要求モジュールは、読み出した保有信号(503)が1(図52の509)か、未確定信号(502)が0(図53の514)ならば処理が終わったので、2サイクル後に確定信号(515)、保有通知信号(図52の516)を1にして処理を完了する。要求モジュールは保有信号(503)が0(図54の517)で未確定信号(502)が1(図54の512)ならば再び要求信号500を1(図54の518)にして送信権を獲得して報告を要求する。モジュールは、未確定信号(502)、保有信号(503)に応答後リング1周のサイクルの2サイクル後の確定信号515が1(図52、図53の515)になったら保有信号(503)、未確定信号(502)の処理を中止或いは完了する。確定信号(504)が0(513)なら、要求モジュールからの再要求に備える。モジュールは、保有信号(503)が1(509)であれば最終状態であることをスヌープできる。メモリモジュールは、読み出しデータがあり、保有通知信号(516)が1(図52の516)であれば送信を抑止する。メモリモジュールは、読み出しデータがあり保有通知信号(516)が0(図53の519)であれば送信する。

【0117】[論理的順序付け] キャッシュの所有権を決めるため、全報知のパケットはその論理的な順序付けを行い、各モジュールが認識する必要がある。全報知のパケットは、リングを1周する間送信権が獲得されているので必ずアドレス0のモジュールを通過する。アドレス0のモジュールはパケットが通過する順に論理的順序を定めることができる。受信モジュールはアドレス0のモジュールの論理的順序に合わせる。受信モジュールとアドレス0のモジュールとの間のモジュールが送信モジュールがあると、アドレス0のモジュールを通過して、受信モジュールを通過し送信モジュールに戻る。パケットの受信時刻で論理的順序を定める。アドレス0のモジュールと受信モジュールとの間に送信モジュールがあると、受信モジュールを通過してアドレス0のモジュールを通過し送信モジュールに戻る。アドレス0のモジュールを通過する順序は、仮想的に送信モジュールから受信モジュールまでパケットを送信する場合の順序になる。パケットの受信時刻にリング1周のサイクル数を加えた

時刻で論理的順序を定める。送信モジュールのアドレスは、リングの宛先のアドレス線に示されている。図55は、アドレス0のモジュール550を通過するパケットA、B、C、D、E、Fが各モジュールにおいて受信される順序(551)を示す。受信モジュールは、自アドレスと宛先アドレスを比較して論理的順序を定める時刻を(553)になるように定める。パケットA(552)はアドレス2のモジュールが送信することを示す。図56は要求元アドレスと自アドレスを比較して自アドレスが大きければ論理的順序づけの時刻は現時刻にリング1周のサイクル数を加えたものになることを示す。

【0118】[選択リング]次に選択リングについて説明する。基本リングの選択線による送信権獲得の代わりに選択リングを用いてアドレス0のモジュールで集中的に送信権の獲得を行うことができる。集中的に送信権獲得を行うために1サイクルの間に2つのパケットの転送を行うことができるのでスループットを向上できる。また、アドレス0のモジュールで処理するので、論理的順序付けは送信権の獲得順でよく、各モジュールで特別の処理はいらない。選択線は分散の送信権獲得を行うのに対し、選択リングは集中的送信権獲得を行う。要求モジュールからアドレス0のモジュールまでの転送時間、アドレス0のモジュールでの送信権処理時間、アドレス0のモジュールから要求モジュールまでの転送時間がかかるのでレイテンシは増加する。信号線の数も増加する。選択線を用いた送信権はリング1周の間、獲得される。全報知のパケット以外は要求元から宛先アドレス迄リングを占有する。宛先アドレスから要求元までは、別のパケットを送信できる可能性がある。選択線を用いた送信権制御は各サイクルの送信権をリング1周のサイクルで処理する。送信権制御を集中するとリングの利用可能の領域に別のパケットを送信できる。選択リングを設け、集中送信権制御をアドレス0のモジュールで行う。集中送信権制御では、1つの送信権の制御にリングの1周のサイクルを使っているのを2つのパケットの送信権制御を1サイクルで処理する。各モジュールからアドレス0のモジュールには各モジュールからの送信要求を報告し、アドレス0のモジュールから獲得した送信権を通知する。キャッシュ制御を行うマルチプロセッサシステムでは、読み出しとキャッシュのライトバックは、同時に発生することが多く、また通常のパケットは全報知である場合が多い。全報知でない2つのブロックパケットをリング1周に並行に転送し通常のパケットはアドレス0から最終アドレスのモジュールの間に1つだけ転送する方式を述べる。論理的順序づけは送信権の順に行われる。ブロックパケットの予約線とそれに付随する信号線は必要でない。各モジュールからアドレス0のモジュールへの送信要求の報告信号は、選択線(1本)及び送信権報告線(10本)である。送信権報告線(10本)には、送信要求報告の時には送信要求(1本)、要求元ア

ドレス(3本)、宛先アドレス(3本)、ブロックパケット(1本)、全報知(1本)を乗せる。送信権報告線(10本)には、ビジー報告の時には送信要求を0にしてビジー報告(1本)、ビジー報告アドレス(3本)、ブロックビジー(1本)、ビジー(1本)を乗せる。送信権報告線には、ビジー解除報告の時には送信要求を0にして、ビジー解除報告(1本)、ビジー解除報告アドレス(3本)、ブロックビジー(1本)、ビジー(1本)を乗せる。アドレス0のモジュールから各モジュールへの送信要求通知信号は、送信権有効(1本)、送信権通知線(12本)である。送信権通知線は、送信要求通知の時にはそのサイクルで獲得した送信権0と送信権1の2つまで通知できる。送信権0は送信要求通知0(1本)、要求元0アドレス(3本)、ブロックパケット0(1本)、全報知0(1本)を乗せる。送信権1は送信要求通知1(1本)、要求元1アドレス(3本)、ブロックパケット1(1本)、全報知1(1本)を乗せる。送信権通知線は、ビジー通知の時には送信要求通知0を0にして1つのモジュールに対するビジー通知(1本)、ビジー通知アドレス(3本)、ブロックビジー(1本)、ビジー(1本)を乗せる。送信権通知線には、ビジー解除通知の時には送信要求通知0を0にして、1つのモジュールに対するビジー解除通知(1本)、ビジー解除通知アドレス(3本)、ブロックビジー(3本)、ビジー(1本)を乗せる。必要な信号線の本数を要約すると報告に11本、通知に13本の合計24本である。2サイクルでパケットを転送する場合には13本が追加になる。フロー制御における送達確認は、1サイクルに2つ並行転送することから受信不能(1本)、受信不能通知(1本)をそれぞれ2本に拡張する。各モジュールからアドレス0のモジュールへパケットの送達確認のため送信権0、送信権1に対応して受信不能報告0(1本)、受信不能報告1(1本)を設ける。各モジュールからアドレス0のモジュールへパケットの送達確認のため送信権0、送信権1に対応して同様に受信不能通知0(1本)、受信不能通知1(1本)を設ける。各モジュールは、送信要求数と処理完了の応答数の差から未処理送信要求数を監視してアドレス0の送信要求バッファが満杯にならぬように送信要求を制御する。同一の要求元からの全報知のパケット、ブロックパケットは要求の順に応答するので送信要求を要求番号まで特定する必要はない。

【0119】[送信権獲得処理]通常のパケットの集中送信権制御の処理は基本リングで行われている処理をアドレス0のモジュール内の論理回路だけで実現する。アドレス0のモジュールでの集中送信権制御を以下に示す。サイクルのあるアドレス(起点)以降からそのサイクルの最後のアドレス(終点)の送信権を定める。選択する送信要求がなければ、アドレス0をリング1周のサイクル後の送信権処理の起点としてそのサイクルの処理



を完了する。選択する送信要求の宛先がリングのアドレスの終点より先にあれば、宛先のアドレスをリング1周のサイクル後の送信権処理の起点として処理を完了する。その宛先が終点より前にあれば、アドレス0をリング1周のサイクル後の送信権処理の起点としてそのサイクルの処理を完了する。

【0120】[高速の送信権獲得回路] 送信権獲得のクリティカルパスは2入力の論理和回路の従属接続にある。図57に各モジュール毎の送信権獲得論理を示す。図58に図57の送信権獲得論理を8個従属接続して各モジュールの送信要求から送信権を獲得する論理を示す。図59に接続数の多い場合に図58の回路を並列に4個接続して送信権を獲得する論理を示す。

【0121】[並列転送ブロックの選別] ブロックパケットに対してはリング1周の間に最大2つのパケットの送信権を処理する。最初のブロックパケットに対してそのサイクルのあるアドレス(起点)以降からそのサイクルの最後のアドレス(終点)の送信権を通常のパケットと同様に定める。送信権を獲得したブロックパケットがあれば、そのブロックの宛先アドレスから要求元迄の間のリングの領域を使用するブロックパケットを並列転送可能なブロックパケットを選別する。図60に送信要求の中で並列転送可能なパケットブロックの選別回路を示す。

【0122】[並列転送ブロックの選択] 並列転送可能なパケットブロックのうちでアドレスの最も小さいブロックパケットが送信権を獲得する。図61は各モジュールでの並列転送可能ブロックから並列転送ブロックを1つ選択する回路を示す。図59の並列転送可能ブロックの選別回路を8個従属接続する。ブロックパケットの送信権獲得のクリティカルパスは2入力の論理和回路の従属接続にある。標準的なゲートによらず専用の回路を用いて良い。

【0123】[並列転送ブロックの通知] 2つのブロックパケットが同一サイクルにあれば、送信権獲得を同時にリングを用いて通知する。2つ目のパケットが次のサイクルにあれば、次のサイクルで通知する。ブロックパケットは少なくとも2サイクル続くので、その送信権処理に2サイクル使用する。

【0124】[フロー制御] フロー制御として、平均ビジー率の低い場合は全報知を優先し、平均ビジー率の高い場合はブロックパケット、全報知でないパケットを優先する。平均ビジー率が高い場合は同のモジュールで定められた数を越えて連続して送信権を獲得できない。

【0125】[階層接続] リング接続したプロセッサシステムの能力を向上させるために階層接続する。プロセッサシステムを副リングとして、副リング間を主リングで階層的に接続する。副リング内で閉じる転送は主リングに影響しない。全体のレイテンシは主リングと副リングのレイテンシの和でレイテンシが決まる。主リングか

らのスループット要求を満たすために副リングのモジュールの接続数を低く押える必要のある時には、副リングをケーブル等で短絡してレイテンシを改善できる。副リングは、多くの場合1つのバックプレーンに收容される。バックプレーン間を接続する主リングは、リングの一つの特徴であるケーブルで接続できるので容易に実装できる。規定の長さを越えて接続する場合は、延長のためのモジュールを追加する。図62に階層接続の接続方法を示す。主リングと副リングの接続機構は、副リングのモジュールのアドレス0を持ち、副リングへの送信権獲得の通知と主リングへの送信要求を同時に行うことができる。接続機構はパケットのヘッダより主リング内のアドレス或いは副リング内のアドレスを生成する。なお、副リングは通常バス構造でも可能である。メッセージのセマフォア、アドレスマッピングレジスタ、DMAレジスタ等IOアドレス空間を共有する場合は、キャッシュしないので宛先が1つに定められパケットの転送は1対1になる。分散共有メモリのように全報知のパケットが生じる場合は、各副リングに必要なスループットは接続する副リングの総数に直線的である。

【0126】図63に階層接続のキャッシュ可能な共有メモリを有する接続機構を示す。副リング700内で閉じるパケットは、主リング701には転送しない。当該副リング700に宛てた主リング701のパケットは、副リング700に転送する。全報知のパケットは、すべての副リング、主リングに転送する。レイテンシを短くできるがリングに必要なスループットは増加する。接続機構にSCIに定められたディレクトリ、プロトコルを用いれば全報知パケットをなくすることができるが、レイテンシは増加し、ディレクトリは大きくなる。リングの全報知能力を活かして接続機構にキャッシュ・コヒーレンスのための全報知パケットを削減し、レイテンシの増加を抑制し、適当な大きさのディレクトリを持つことができる。副リング700に属するモジュールがその副リングに属するローカル・メモリ703にアクセスする場合、他の副リング704のキャッシュに写しがなければ主リング701にそのパケットを転送しないでよい。副リング700に属するローカルメモリ703の写しを他の副リング704のキャッシュに持っていないことを示すローカル・ディレクトリ705を持てば、主リング701へのパケットの転送を削減できる。ある副リング700に属するモジュールが、他の副リング706に属するローカル・メモリ707にアクセスする場合、アクセスされたローカルメモリ707を持つ他の副リング706がその副リング704には写しを持っていないことを通知すれば、その副リング704内にパケットを転送しないでよい。他の副リング706に属するローカルメモリ707が、その副リング704のキャッシュに写しのある状態を示すリモート・ディレクトリ706を持てば、その副リング704内に転送するパケットを削減で

きる。宛先の副リングは36ビットの内の最初の4ビットを用いて行われる。キャッシュ可能なメモリ領域であり実装されていることを確認してからディレクトリにアクセスする。

【0127】[ローカル・ディレクトリ] メモリアドレスは36ビット(64GB)で、副リングに実装するローカルメモリのアドレスは32ビット(4GB)、ラインの大きさは6ビット(64B)とする。ローカル・ディレクトリのエントリ数は、32ビットから6ビットを引いた26ビット(64Mエントリ)である。

【0128】図64にアドレスの関係を示す。ローカル・ディレクトリのエントリには、写しビットと変更ビットの2ビットを持つ。写しビットは、他の副リングのキャッシュに写しがあれば1、なければ0を格納する。変更ビットは、他の副リングのキャッシュに写しがあり変更されている可能性があれば1、なければ0を格納する。ローカル・ディレクトリは、ローカルメモリの256分の1の容量である。

【0129】図65にローカル・ディレクトリを示す。メモリアドレスを40ビット(1024GB)に、副リングに実装するローカルメモリのアドレスを36ビット(64GB)に拡張する場合は、ローカル・ディレクトリのエントリ数は36ビットから6ビットを引いた30ビット(1024Mエントリ)になる。大容量であるので、ローカル・ディレクトリ・テーブルとしてローカルメモリ内に格納してローカル・ディレクトリにはその写しを持つ。ローカル・ディレクトリの写しは、ラインの大きさと同じ64B単位の大きさにする。64B(512ビット)には256(8ビット)のエントリを収容できる。ローカル・ディレクトリには、26ビットから8ビットを引いた18ビット(256Kエントリ)の写しを格納できる。ローカルメモリのアドレスの36ビットから8ビットと6ビットを引いた22ビットを18ビットに変更するためにアドレス変換テーブルを持つ。

【0130】図66にアドレスを拡張する場合を示す。アドレス変換テーブルは、256Kエントリ(18ビット)を持ち、各エントリには22ビットのアドレスを持つ。アドレス変換テーブルの1つのエントリは、22ビットから18ビットを引いた4ビットに相当する16単位の内の1つの写しを示す。

【0131】図67にアドレス変換テーブルを示す。ローカルディレクトリは、ローカルメモリのすべての範囲のキャッシュの状態を格納しているのでディレクトリリングは必要ない。アドレス変換テーブルのエントリ数も、18ビットのアドレスをハッシュすることにより減らせる。

【0132】[リモート・ディレクトリ] 各エントリには、副リングに属するローカルメモリの写しがその他の副リングにあれば1、なければ0を主リング内の副リングのアドレス0からアドレス7まで順に格納する。この

8ビットを写しベクトルと呼ぶ。主リング内の副リング(アドレス)を31に拡張するときには、各ビットをグループ単位に割り当てる。グループ内はSMPでグループ間はクラスタの場合等に効果がある。例えば第1ビット目は、アドレス0-3のグループの様に割り当てる。4つのアドレスの中の1つにしか写しがなくとも、4つのアドレスに対応する副リングすべてに写しがあるとして取り扱う。リモートディレクトリのエントリ数は、再利用する頻度によって決められる。22ビット(4Mエントリ)にする場合は、そのエントリを特定するためにローカルアドレスの32ビットから22ビットと6ビットを引いた4ビットを各エントリに格納する。ローカルアドレスを36ビットに拡張するには、エントリを特定するために36ビットから22ビットと6ビットを引いた8ビットを格納する。図68にリモート・ディレクトリを示す。

【0133】[写しベクトルの通知] 副リングに属するローカルメモリにアクセスする全報知のパケットを主リングに転送する時には写しベクトルをヘッダに格納して転送する。副リングに属するローカルメモリにアクセスする全報知のパケットを受け取った時には写しベクトルをパケットを用いて通知する。写しベクトルは自身のアドレス位置には写しビットはいらないのでその位置が1であれば有効であることを示す。

【0134】[ローカルメモリの種別] ローカル・ディレクトリによって、ローカルメモリは写し無し(N)、写し有り(C)、変更有り(M)に分けられる。リモート・ディレクトリによって、ローカルメモリの写し有り(C)の一部は写しベクトル有り(CH)、変更有り(M)の一部は写しベクトル有り(MH)に分けられる。図69にディレクトリにより分けられるローカルメモリの種別を示す。

【0135】[ローカルメモリアクセス] 副リング705に属するモジュールがローカルメモリ703を写すパケットを送信して、接続機構708がこれを受信したら、N、Cなら主リングへ転送しない。Mなら変更した内容をキャッシュより書き戻す必要がある。このため、写しベクトルをヘッダに格納して主リング701に全報知のパケットを送信する。主リング701の他の副リング704は、写しベクトルが有効で且つ0であれば副リング704内にパケットを転送しない。写しベクトルが有効で且つ0が無効であれば副リング704内にパケットを転送する。他の副リングから書き戻される間、接続機構において同一アドレスに対するアクセスは待たされる。書き戻されるとCHとなり処理が引き継がれる。副リング700に属するモジュールがローカルメモリ703を変更するパケットを送信して、接続機構708がこれを受信した場合、Nなら主リングに転送しない。Cならキャッシュの写しを廃棄する必要がある。Mならキャッシュの内容を書き戻した上で廃棄する必要がある。

C、Mなら写しベクトルをヘッダに格納して主リング701に全報知の packets を送信する。主リング701の他の副リング704は、写しベクトルが有効で且つ0であれば副リング704内に packets を転送しない。写しベクトルが有効で且つ0が無効であれば副リング704内に packets を転送する。Mなら書き戻されるのを待ちNに、CならそのままNになる。以上で明らかな様に、ローカルメモリで閉じる処理は他の副リングの影響を受けない。

【0136】 [リモートメモリ・アクセス] 副リング700に属するモジュールが他の副リング706に属するローカルメモリ707 (リモートメモリ) を写す packets を送信して接続機構708がこれを受信したら、主リング701に転送する。ローカルメモリ707を有する副リング706は、写しベクトルを packets により通知する。主リング701の他の副リング704は、写しベクトルの受信してから受信してある packets の処理を行う。写しベクトルが有効で且つ0であれば副リング704内に packets を転送しない。写しベクトルが有効で且つ0が無効であれば副リング704内に packets を転送する。ローカルメモリ707の状態はCHになる。副リング700に属するモジュールが他の副リング706に属するローカルメモリ707 (リモートメモリ) を変更する packets を送信して接続機構708がこれを受信したら、主リング701に転送する。ローカルメモリ707を有する副リング706の接続機構709は、 packets を受信すれば副リング706内に転送する。ローカルメモリ707を有する副リング706の接続機構709は、Nならローカルメモリ707をアクセスし、また他の副リング704で packets の処理をしないで良いことを packets により通知する。Cならローカルメモリ707をアクセスし、また他の副リング704のキャッシュに写しがあれば廃棄することを通知する。Mならば、書き戻しを行ってから廃棄することを通知する。これらのため、キャッシュの写しベクトルを packets により通知する。主リング701の他の副リング704は、 packets から写しベクトルを受信してから受信してある packets を処理する。写しベクトルが有効で且つ0であれば副リング704内に packets を転送しない。写しベクトルが有効で且つ0が無効であれば、副リング704内に packets を転送する。ローカルメモリの種別はMHになる。以上で明らかな様にリモートメモリにアクセスする処理は、レイテンシを遅くすることなくまた必要のない副リングへのアクセスを行わない。

【0137】 [双方向リング] モジュール間は、1対1の1方向の伝送を行っている。双方向伝送回路を用いて2つの独立する双方向リングを構成することができる。双方向リングはリング接続の配線をそのままにして性能の限界を2倍にする増設性を与える。2つのリングの同期ができないと言う制約があり、伝送の方向性の弁別回

路の弁別精度、周波数特性の制約から動作周波数の限界は低くなる。また弁別回路そのもの及び弁別回路が発生する雑音のために中継モードのモジュールの数の制限が厳しくなり、転送モードのモジュールの割合が増えてレイテンシが長くなる。全報知の通常 packets は1つのリングを利用して論理的順序性の保障を行い、順序性の制約のない通常の packets 、ブロック packets を2つのリングに分散させる。2つの並列接続リングの構成に比較すれば、性能は若干低くなる。図70に弁別回路の例を示す。図71に双方向リングのモジュール間の接続を示す。

【0138】 [構成変更] リング制御ユニットは、各モジュールへのヘルスチェックのフレームを送信して、応答のないことで構成変更が行われていることを検知する。リング制御ユニットは、各モジュールへのヘルスチェックを行い、再応答により構成変更が完了したことを知り、イニシャライズを行う。構成変更によって一般にいくつかの packets が失われるので、適当なリスタート・ポイントまで遡る必要がある。リングの構成変更の際に、計画停止してからリングの動作を停止しイニシャライズを行った後に直ちに再開する方法を図72に示す。モジュールの引き抜きを早期に知るために、モジュールの引抜レバーの引き起こしを引抜センサで検知してモジュール制御ユニットを通じてリング制御ユニットに報告する。リング制御ユニットは、引き抜きの報告により各モジュールに対して新たな動作要求を開始せず、それ以前の動作要求への応答は行いしかる後ハードウェア動作を正常停止することを指示する。モジュールの引抜に連動して短絡スイッチは短絡する。リング制御モジュールは、ヘルスチェックによりモジュールの引抜の完了を自動的に確認する。モジュールの引抜中もリングのクロックは供給されるので引抜モジュール以外のモジュールの正常停止状態は継続する。モジュールの引抜の報告よりアドレスを再設定し、スキュー補正の必要なモジュールに対してイニシャライズを行い、回送モードのモジュールでリングの1周のサイクル数を決める。リング制御ユニットは、ハードウェア停止状態を解除し、新たな構成情報を伝える。モジュールの挿入に際してもリングの動作は停止する。モジュールの挿入を早期に知るために挿入信号 (1本) を設ける。挿入信号線はそのモジュールのガイドから始まり、短絡コネクタがあればこれを經由して次のモジュールに接続して終わる。挿入信号は、モジュール間の接続を行うのでモジュールは転送しない。モジュールのガイドを通過することを挿入センサで検知して、挿入信号を一時的に短絡する。各モジュールは、挿入信号が短絡されるかを常時監視していて、短絡されれば前のモジュールが挿入されるはずであるということを経由して次のモジュールに報告する。リング制御ユニットは、各モジュールに対して新たな動作要求を開始せず、以前の動作要求への応答を行いしかる後ハードウェア動

作を正常停止することを指示する。モジュールの挿入に連動して、短絡スイッチを機械的にオフにする。リング制御ユニットは、ヘルスチェックの応答により挿入の完了とそのモジュールを特定する。モジュールの挿入中も、リングのクロックは供給されるので挿入モジュール以外のモジュールの正常停止状態は継続する。挿入の報告よりアドレスを再設定しスキュー補正の必要なモジュールに対してイニシャライズを行い、回送モードのモジュールでリングの1周のサイクル数を決める。リング制御ユニットは、ハードウェア停止状態を解除し新たな構成情報を伝える。2つのリングを接続して、通常状態では同期、並行動作を行い、故障時は縮退運転を行うことができる。2つのリングの同期のために回送モードにおいて1サイクルの遅延ラッチを設けリングの1周のサイクル数を一致させ、リング制御ユニットの動作を同期させる。階層接続の主リングを複数持てば拡張性は若干失われるがスループット、信頼性が向上する。共通の基板によらず、各スロット間をケーブルで接続することができる。

【0139】[インターフェイス信号線の具体例] インターフェイス信号の具体例を図73に示す。この例では、メモリにディレクトリを持つのでディレクトリ・リングはない。また、選択リングは採用していない。

【0140】[イニシャライズ] 各モジュールのアドレス、回送モード、転送モード、中継モード、位相調整回路、サイクル調整回路、可変遅延回路、リングを一周するサイクル数、複数サイクルのパケットのクロックの同期は、イニシャライズ時に設定する。リング制御ユニットは、リング全体、モジュール間の相互関係、処理の同期等を行うため各モジュールのモジュール制御ユニットに指示を出す。モジュール制御ユニットは指示に基づき可変遅延回路の設定などの処理を行う。

【0141】[リング制御ユニット] リング制御ユニットは、リング制御のイニシャライズを制御、統括するとともに、クロックを各モジュールにリング状に配信する。リング制御ユニットはクロックが配信されていることを確認するために受信する。各モジュールは、リング制御ユニットからの指示を基に処理し、リング制御ユニットに状態を応答する。リング制御ユニットは、フレームを用いて指示を伝達する。指示は、ステップバイステップで行う。フレームの内容は、コマンドバイト、アドレスバイト、2バイトのギャップ、2バイトのデータである。コマンドはリング制御ユニットの指示、応答要求である。アドレスは各モジュールあるいは全報知を指定する。2バイトのギャップは、モジュール制御ユニットに応答の処理する時間を与える。データは、各モジュールへの指示或いは各モジュールからの応答である。指示のフレームではアドレスで指定されたモジュールはフレームのデータを受信する。応答要求のフレームでは、アドレスモジュールで指定されたモジュールはそのフレー

ムのデータ部に内容を入れて応答する。リング制御ユニットと各モジュールはリング状に接続する。制御リングは、各モジュールとの独立性を高めるためいわゆるリングバスの構造を取っている。フレームは、ビット、バイトからなる。ビット同期は、制御クロックで行う。制御クロックの位相はフレームのビットの中央で立ち上がる。バイトの同期は、同期パターン“000000000000FF”で行う。フレームの同期は、“FExxxxxx”で行う。フレーム間は“00”を埋める。

【0142】[パワーオン・モード] リング制御ユニットは、全報知で各モジュールにパワーオン・モードを指示する。パワーオンモードでは回送モードと転送モードのみを使用して動作の変動に対処する。

【0143】[モジュールアドレスの設定] リング制御ユニットは、各モジュールに全報知でアドレス設定準備の指示を出す。各モジュールは、あらかじめ挿入するスロット番号をもっている。アドレス設定準備の指示を出されると各モジュールはアドレス部に対してスロット番号で動作する。リング制御ユニットは、アドレス部にスロット番号を入れてモジュールの種類、論理アドレス等の応答を要求する。指定されたモジュールは、データ部にプロセッサ、I/O、メモリ、プロセッササブシステム等の種類を応答する。リング制御ユニットは、アドレス部にスロット番号を入れてデータ部にモジュールのアドレスを入れアドレス設定を指示する。リング制御ユニットは、各モジュールに全報知でアドレス設定完了の指示を出す。アドレス設定完了の指示を出されると各モジュールはアドレス部に対して設定されたアドレスで動作する。リング制御ユニットは、回送モード、転送モード、中継モードを設定する。

【0144】[受信線毎のスキュー補正] 受信線毎のスキュー補正を行う場合は、リング制御ユニットは回送モード、転送モードのモジュールより1000の繰り返しパターンの送信を指示する。引き続くモジュールに対してスキュー補正を基本動作に分解して指示する。最初に調整クロックの位相を初期設定する。可変遅延回路の遅延量の0設定、位相調整回路の調整開始の指示、受信線のすべての0の検出指示、位相調整回路の遅延量の1ステップ毎の走査、0検出の結果のステップ毎の報告により調整クロックは初期設定される。次に最大の遅延の受信線に遅延クロックの位相を合わせる。受信線のすべての1検出指示、位相調整回路の遅延量の1ステップ毎の走査、1検出結果のステップ毎の報告により位相調整回路の位相が最大の遅延の受信線に合わせられる。次に可変遅延回路の遅延量の1ステップ毎の走査により1検出の最後の遅延量が格納される。遅延量の走査の完了後、格納された遅延量を戻してスキュー補正が完了する。この場合は、遅延クロックの逆位相のクロックの供給により受信位相も同時に決定する。この動作をリングの全モ

ジュールについて行い、スキュー補正の完了を指示して遅延クロックを書くモジュールのリング制御部以外にも供給する。

【0145】[受信位相の決定] 前のモジュールに1000の設定パターンの送信を指示する。最初に最小の遅延の受信線に、遅延クロックを合わせこれを前縁とする。次に最大の遅延の受信線に、遅延クロックを合わせて後縁とする。引き続き受信線の間には遅延クロックを合わせて、その逆位相のクロックを供給して受信位相が決定する。モジュール制御ユニットは、基本動作に分解してモジュール制御ユニットに指示する。

【0146】[サイクル調整回路] 回送モードのモジュールに10の繰り返しパターンの送信を指示する。回送モードのモジュールに、サイクル調整回路に内蔵しているクロックの規定時間前とクロックの規定時間後の入力の検出回路により、それぞれを値を応答させる。2つの値が同じであれば、サイクル調整ラッチを経由せずに回送ラッチにラッチするように指示する。2つの値が同じでなければ、サイクル調整ラッチを経由して回送ラッチにラッチするように指示する。

【0147】[サイクル数の設定] リング制御モジュールは、回送モードのモジュールに1つの1と引き続く指定する数の0の繰り返しパターンの送信を指示する。リング制御モジュールは回送モードのモジュールにリング1周のサイクル数を1を受信する間隔で行うように指示して、その報告を要求する。リング制御ユニットは、全報知で各モジュールにリング1周のサイクル数を知らせる。

【0148】[複数サイクルの位相の設定] 全報知で複数サイクルの位相の設定の指示を出す。サイクル数に合わせた1に続く0の繰り返しパターンを回送モードより送信する。全報知で設定の完了を指示する。

【0149】[動作モードへの移行] リング制御ユニットは、パワーオン直後のスキュー量が安定状態になる規定の時間後にパワーオンモードを解除して動作モードに移行する。リングの各モジュールに対して回送モード、転送モード、中継モード、受信線毎のスキュー補正、受信位相、サイクル調整回路、サイクル数を再度決定する。

【0150】[誤動作検出] 制御リングは、ヘルスチェックにより各モジュールを監視する。基本リング内の誤動作の検出は、データ線についてはエラー検出コード、選択線については送信権放棄が定期的に行われているかの監視で行う。リングの回復不能或いは規定回数以上の誤動作時はモジュール制御ユニットによりリング制御インターフェースを通じてリング制御ユニットに知らされる。リング制御ユニットよりリングの停止を指示してイニシャライズを行う。リング制御ユニットは各モジュールの固有の障害の報告をヘルスチェックにより受けて、モジュール制御ユニットに切り離しを指示する。

【0151】

【発明の効果】本発明のリング接続を用いたデータ転送方法及び情報処理システムによれば、プロセッサを経由するためのレイテンシを短縮することができ、高性能な情報処理システムを実現することができる。

【0152】また、リングが通常動作しなくなる切り替え時にもリングよりモジュールに安定したクロックを供給することができる。

【0153】さらに、本発明のリング接続を用いたデータ転送方法及び情報処理システムによれば、情報処理システムのスループットを向上させることができる。

【0154】さらに、本発明のリング接続を用いたデータ転送方法及び情報処理システムによれば、リングの構成変更に伴うシステム停止時間を短縮することができる。

【図面の簡単な説明】

【図1】リング、バス及びスイッチの各プロセッサ間接続技術の比較を示す図である。

【図2】リング接続におけるリングのインターフェイス部の構成を示す図である。

【図3】転送のビット幅が16ビットのリングを示す図である。

【図4】転送のビット幅が64ビットのリングを示す図である。

【図5】従来技術のリングのフラグ処理を示す図である。

【図6】本発明のリングのフラグ処理を示す図である。

【図7】本発明のリングの送受信を示す図である。

【図8】従来技術のリングのバッファを示す図である。

【図9】従来技術のリングのクロックの乗り換えを示す図である。

【図10】本発明のリングの送受信のクロックを示す図である。

【図11】本発明のクロックの位相の進行を示す図である。

【図12】本発明のクロックの系統を示す図である。

【図13】従来技術のリングの転送を示す図である。

【図14】本発明のリングの転送を示す図である。

【図15】本発明の各モジュールの動作モードを説明するためのパケットの転送例を示す図である。

【図16】本発明のスキュー補正を示す図である。

【図17】本発明のリングの種類を示す図である。

【図18】本発明のブロックパケットリングを示す図である。

【図19】本発明の並列転送の例を示す図である。

【図20】本発明の並列転送パケットの選択論理方式を示す図である。

【図21】本発明のディレクトリリングを示す図である。

【図22】本発明のディレクトリへのアクセスを示す図

である。

【図 23】従来技術の冗長リングを示す図である。

【図 24】本発明の基本リングと制御リングを示す図である。

【図 25】本発明によるレイテンシの改善内容を示す図である。

【図 26】本発明による動作周波数に改善内容を示す図である。

【図 27】本発明による増設性を示す図である。

【図 28】本発明の基本リングを示す図である。

【図 29】本発明の制御リングを示す。

【図 30】従来の送信権獲得方式を示す。

【図 31】本発明のゲート 2 段で転送する方式を示す。

【図 32】従来の受信方式を示す図である。

【図 33】本発明の送受信を並行する方式を示す図である。

【図 34】本発明の選択線、アドレス線とデータ線との関係を示す図である。

【図 35】本発明の受信方式を示す図である。

【図 36】本発明の受信動作のタイムチャートを示す図である。

【図 37】本発明の送信動作の概念図を示す図である。

【図 38】本発明の送信動作を示す図である。

【図 39】本発明の送信動作の一例を示す図である。

【図 40】スキュー補正しない例を示す図である。

【図 41】本発明のクロックのスキュー補正する例を示す図である。

【図 42】本発明の受信線毎にスキュー補正する例を示す図である。

【図 43】従来方式の送受信方式を示す図である。

【図 44】本発明のリングのクロックに同期して動作する方式を示す図である。

【図 45】本発明の転送モードの送信動作の概念図を示す図である。

【図 46】本発明のゲートの遅延を利用した位相調整回路を示す図である。

【図 47】本発明のサイクル調整回路を示す図である。

【図 48】本発明の受信線毎にスキュー補正する転送モードの概念図を示す。

【図 49】本発明の送達確認の説明を示す図である。

【図 50】本発明の独立した送信権が 2 つのブロック転送を示す図である。

【図 51】本発明のリング 1 周より大きなブロック転送を示す図である。

【図 52】本発明のディレクトリリングのキャッシュがある場合の動作の例を示す図である。

【図 53】本発明のディレクトリリングのキャッシュのない場合の動作の例を示す図である。

【図 54】本発明のディレクトリリングの未確定で再要求の動作例を示す図である。

【図 55】本発明の各モジュールにおいて受信される順序を示す図である。

【図 56】本発明の論理的順序の時刻の決定法を示す図である。

【図 57】本発明の送信権獲得論理を示す図である。

【図 58】本発明の 8 モジュールの送信権獲得論理を示す図である。

10 【図 59】本発明の接続数の多い送信権処理を示す図である。

【図 60】本発明の並列転送可能ブロックの選別方式を示す図である。

【図 61】本発明の並列転送ブロックの選択を示す図である。

【図 62】本発明の階層接続の接続方法を示す図である。

【図 63】本発明の階層接続の場合のキャッシュ可能なアクセスをする接続機構を示す図である。

20 【図 64】本発明の階層接続のアドレスの関係を示す図である。

【図 65】本発明のローカルディレクトリを示す図である。

【図 66】本発明のアドレスの拡張を示す図である。

【図 67】本発明のアドレス変換テーブルを示す図である。

【図 68】本発明のリモートディレクトリを示す図である。

30 【図 69】本発明におけるディレクトリにより分けられるローカルメモリの種別を示す図である。

【図 70】本発明の双方向リングの弁別回路を示す図である。

【図 71】本発明の双方向リングのモジュール間接続をしめす図である。

【図 72】本発明のモジュールの挿入、引き抜きの早期検出方式を示す図である。

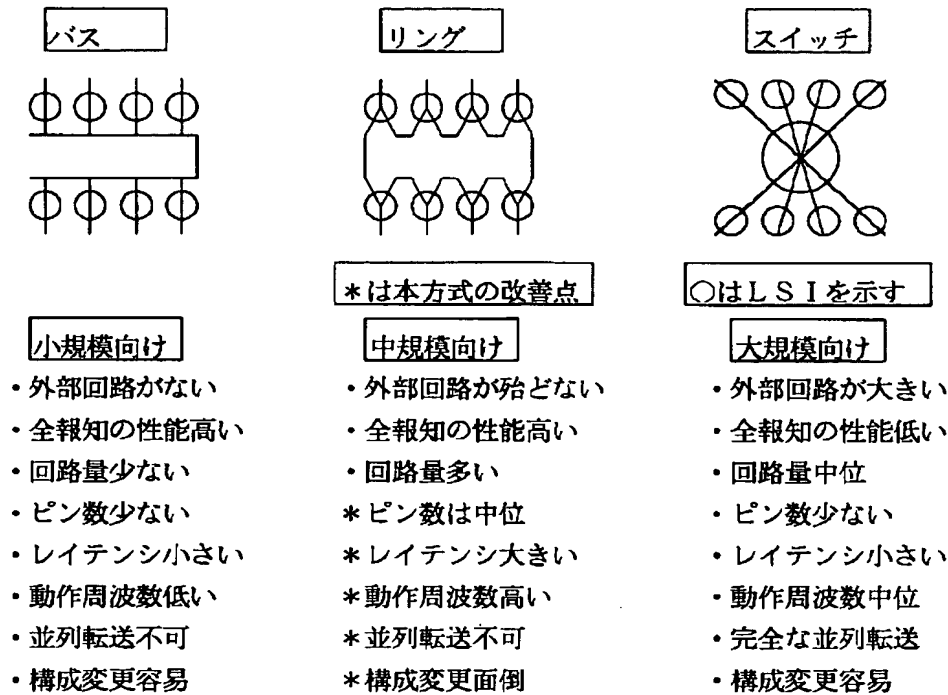
【図 73】本発明のリングインターフェースの具体例を示す図である。

【符号の説明】

40 11…データ線、12、21、22…論理積、13…論理和回路、14、15…受信ラッチ、16…自分アドレス、17…一致論理、18、19、20…ラッチ、23…受信パケット、101、102、104、105…モジュール、701…主リング、700、704、706…副リング、703、707…ローカルメモリ、705、706…リモートディレクトリ、708、709…接続機構、720…選択線、728、736…論理和、729、737…論理積、733…受信線。

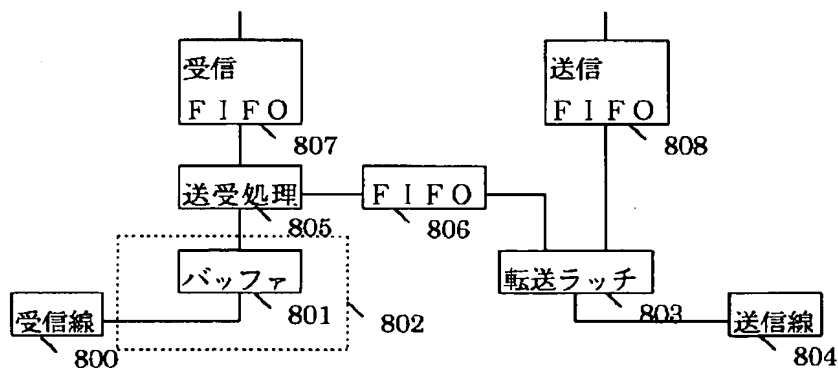
【図 1】

図 1



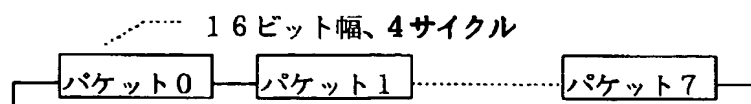
【図 2】

図 2

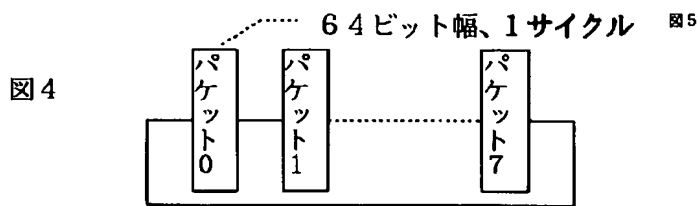


【図 3】

図 3



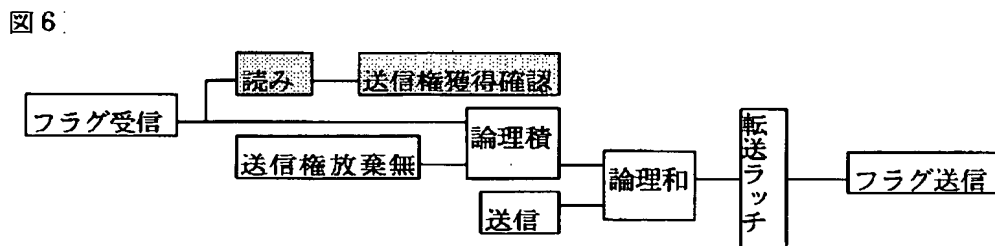
【図 4】



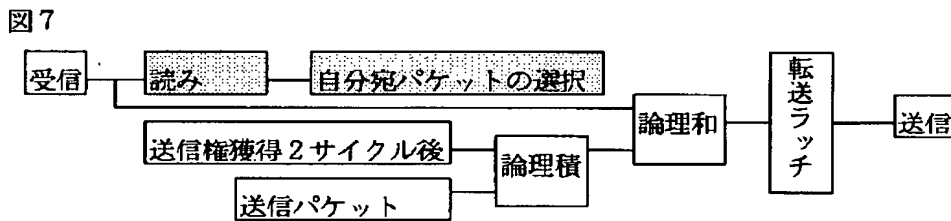
【図 5】



【図 6】



【図 7】



【図 8】

【図 18】

図 8

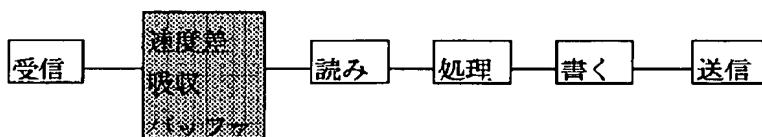
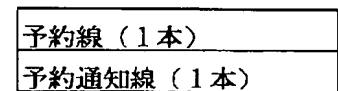
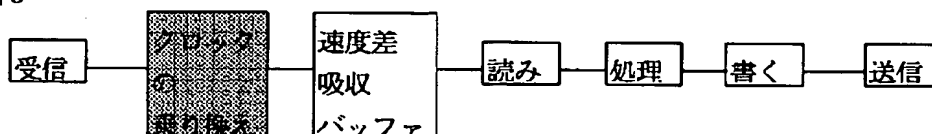


図 18



【図 9】

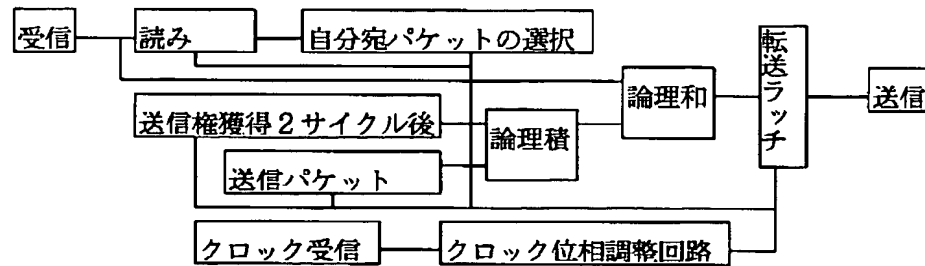
図 9





【図 1 0】

図 1 0



【図 1 1】

【図 1 7】

図 1 1

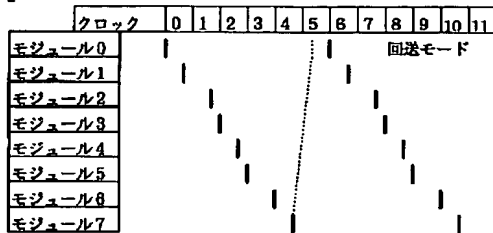
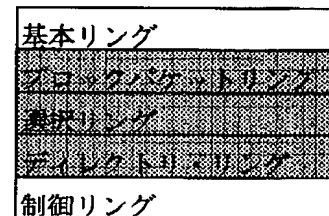
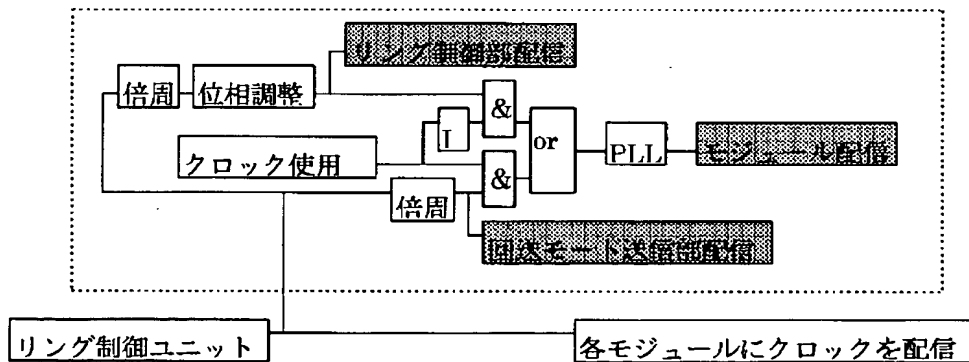


図 1 7



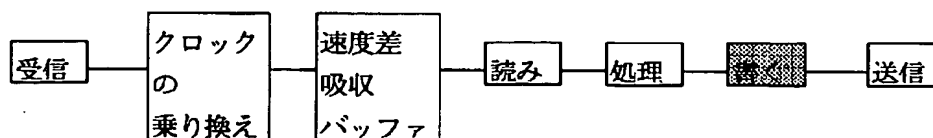
【図 1 2】

図 1 2



【図 1 3】

図 1 3



【図 14】

【図 65】

図 14

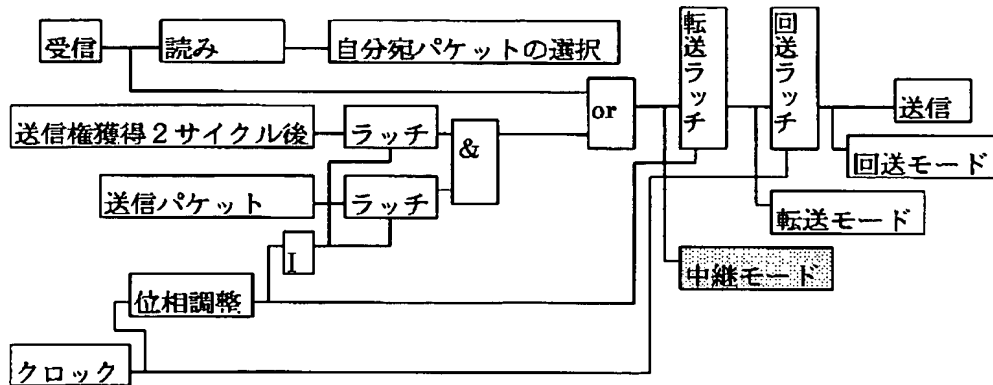
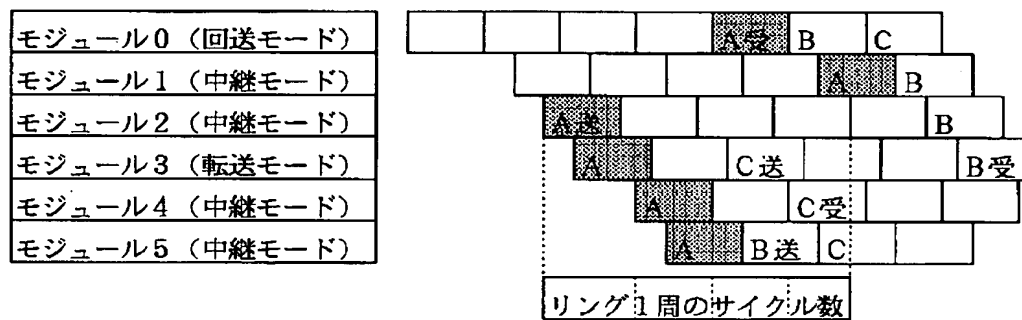


図 65

1	写しビット	変更ビット
2	写しビット	変更ビット
64M	写しビット	変更ビット

【図 15】

図 15



【図 16】

【図 67】

図 16

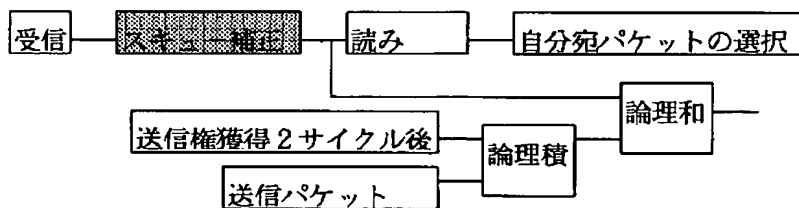


図 67

1	22ビット
2	22ビット
256K	22ビット

【図 23】

【図 64】

図 23

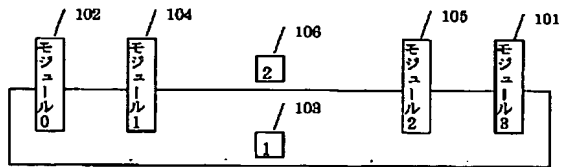


図 64

メモリアドレス (36ビット)
図リングアドレス (4ビット)
ローカル・メモリアドレス (32ビット)
ローカル・ディレクトリ (26ビット)

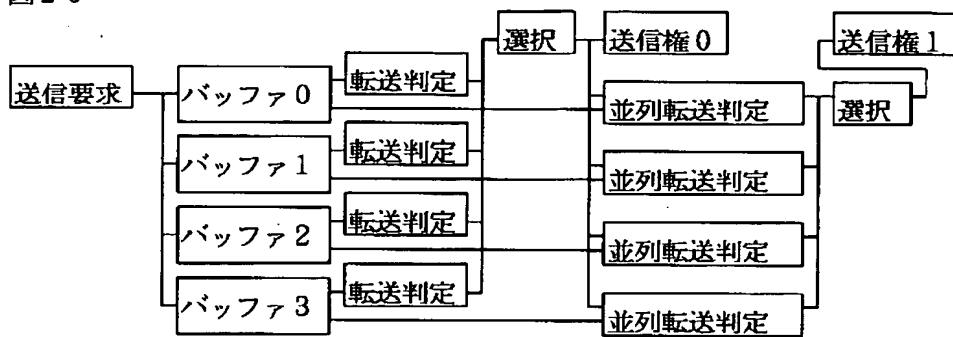
【図 1 9】

図 1 9



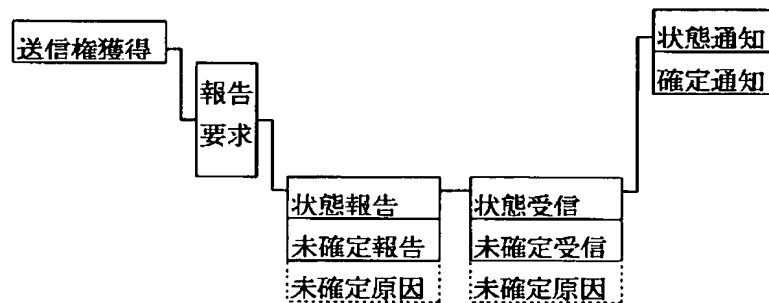
【図 2 0】

図 2 0



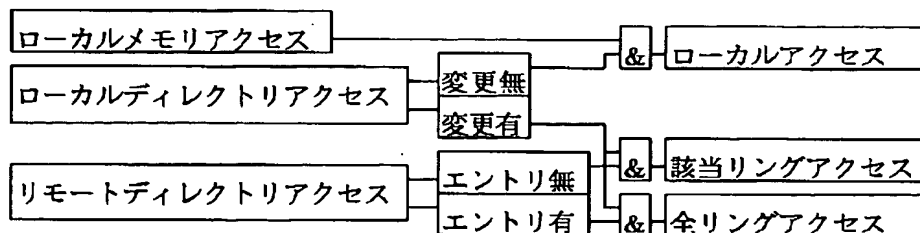
【図 2 1】

図 2 1



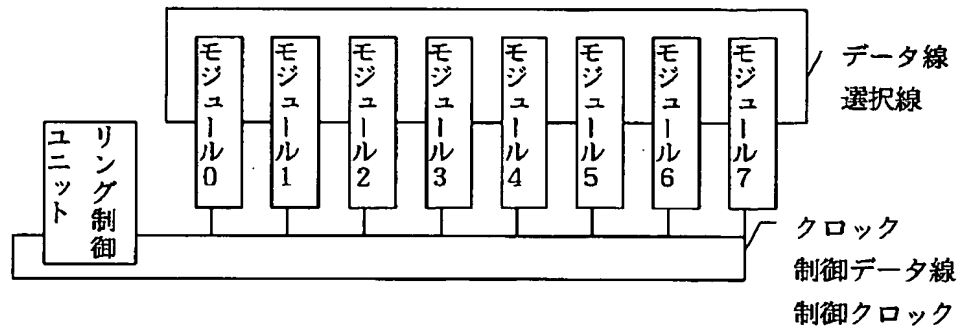
【図 2 2】

図 2 2



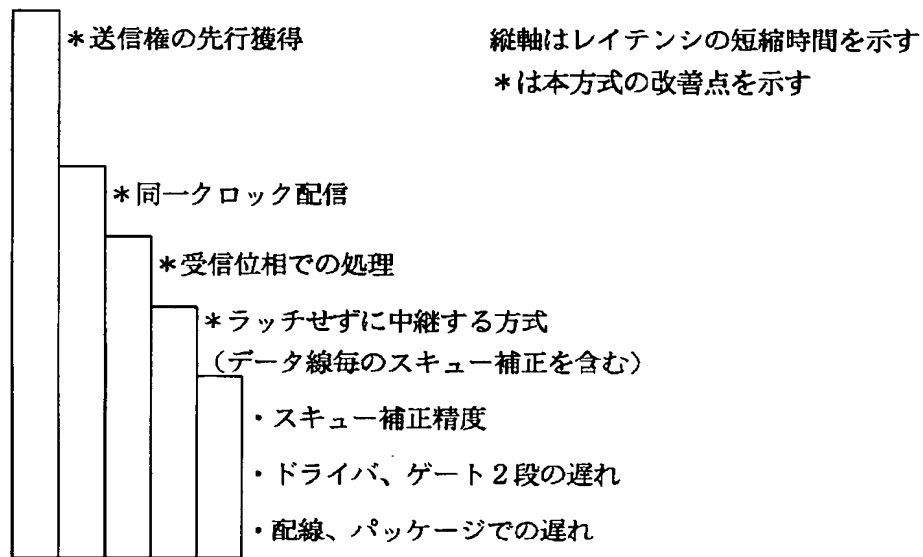
【図 2 4】

図 2 4



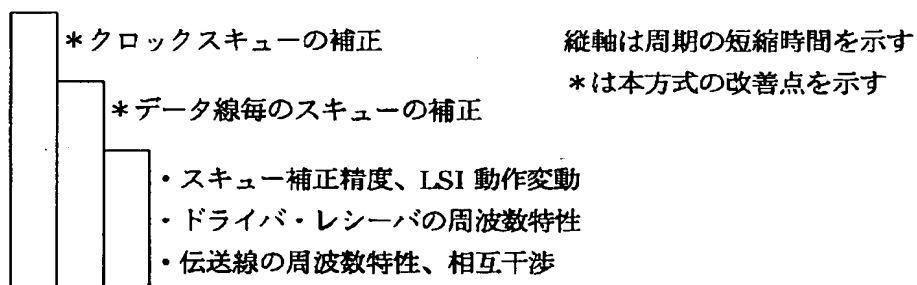
【図 2 5】

図 2 5



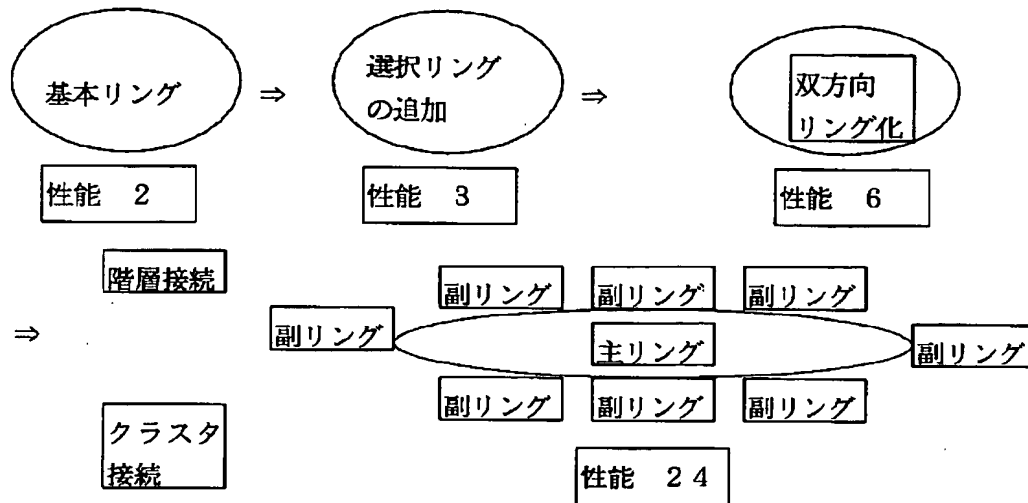
【図 2 6】

図 2 6



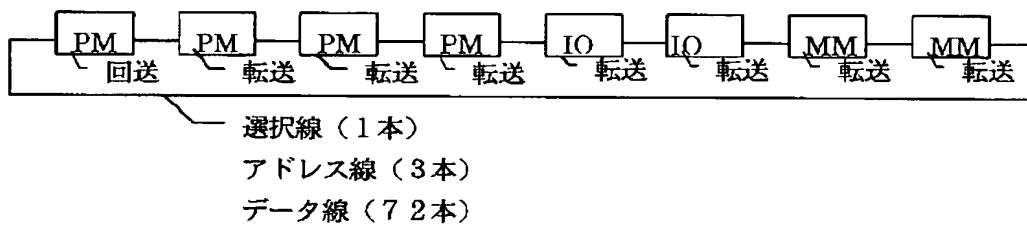
【図 2 7】

図 2 7



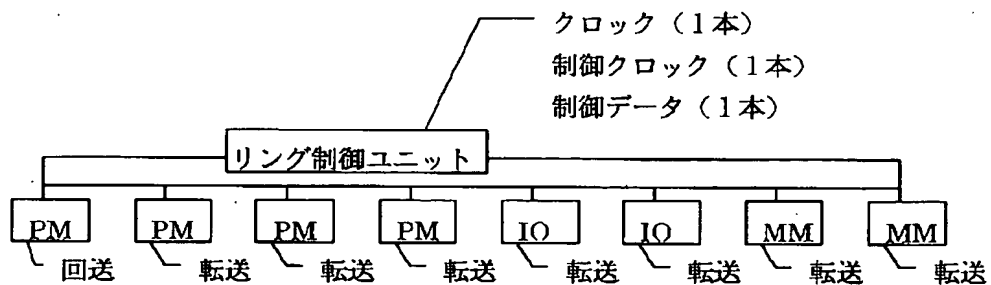
【図 2 8】

図 2 8



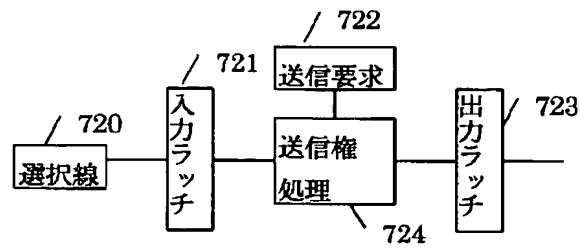
【図 2 9】

図 2 9



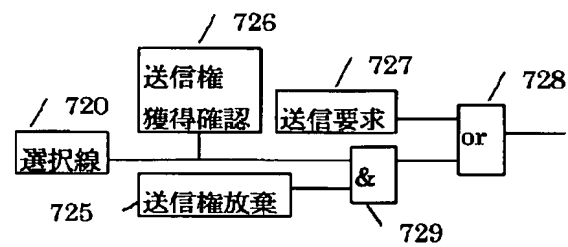
【図 3 0】

図 3 0



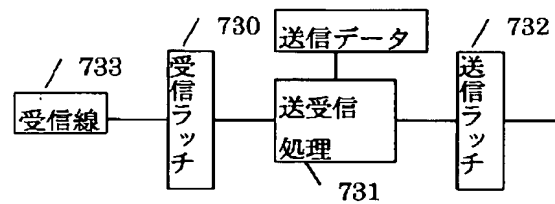
【図 3 1】

図 3 1



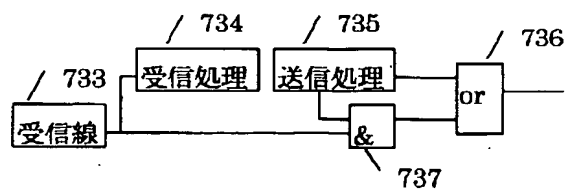
【図 3 2】

図 3 2



【図 3 3】

図 3 3



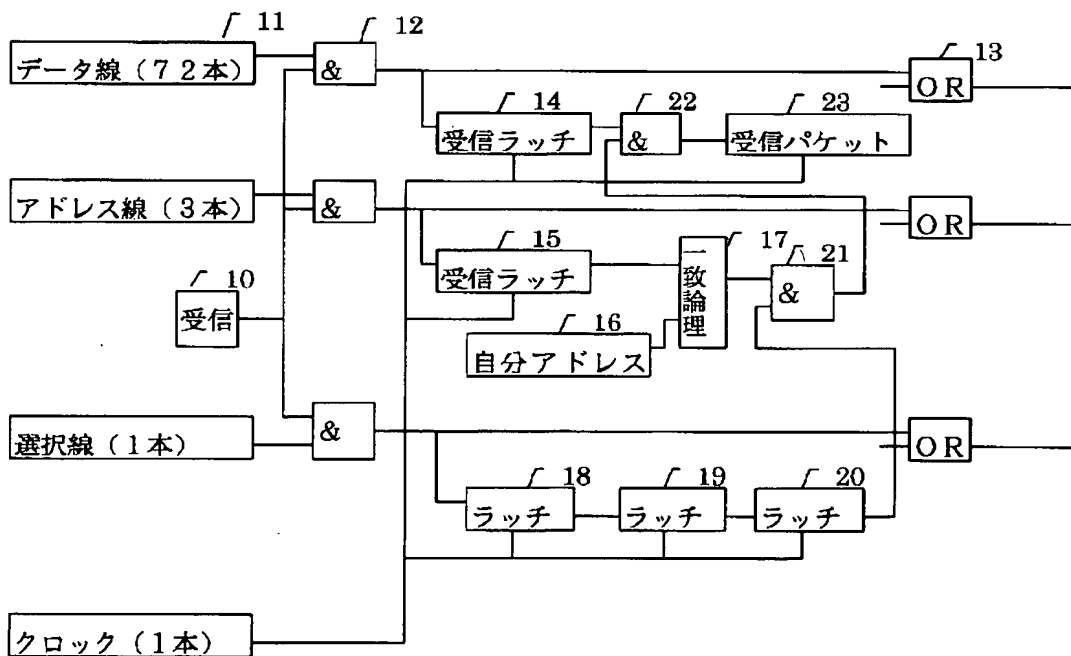
【図 3 4】

図 3 4

時間	クロック	選択線	アドレス線	データ線
1	1	0	x	x
2	1	1	x	x
3	1	0	x	x
4	1	1	3	A
5	1	1	x	x
6	1	0	5	B

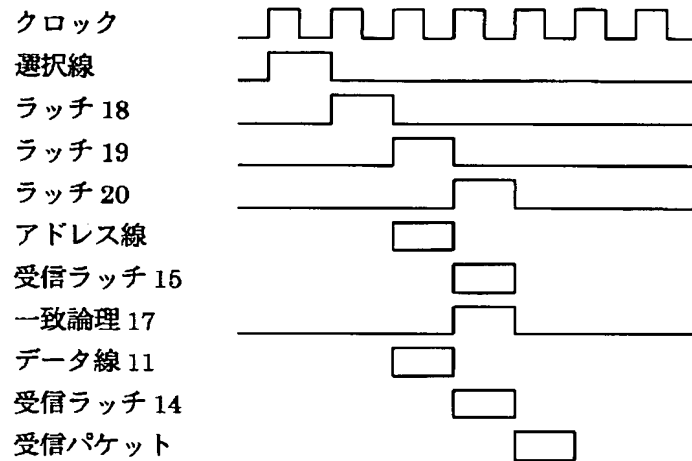
【図 3 5】

図 3 5



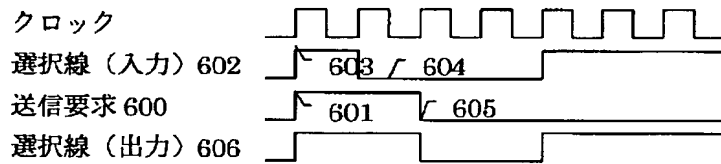
【図 3 6】

図 3 6



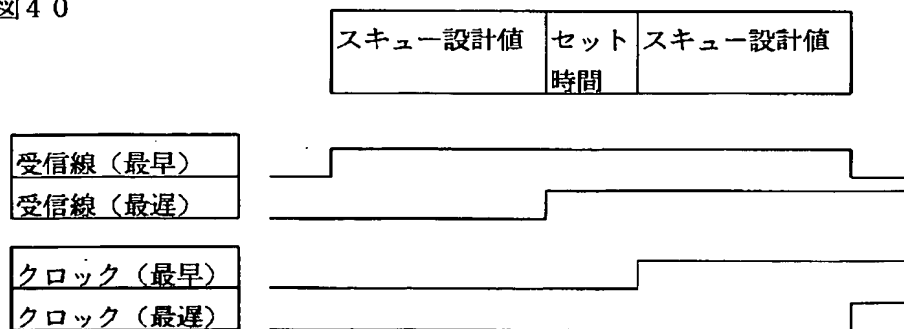
【図 3 7】

図 3 7



【図 4 0】

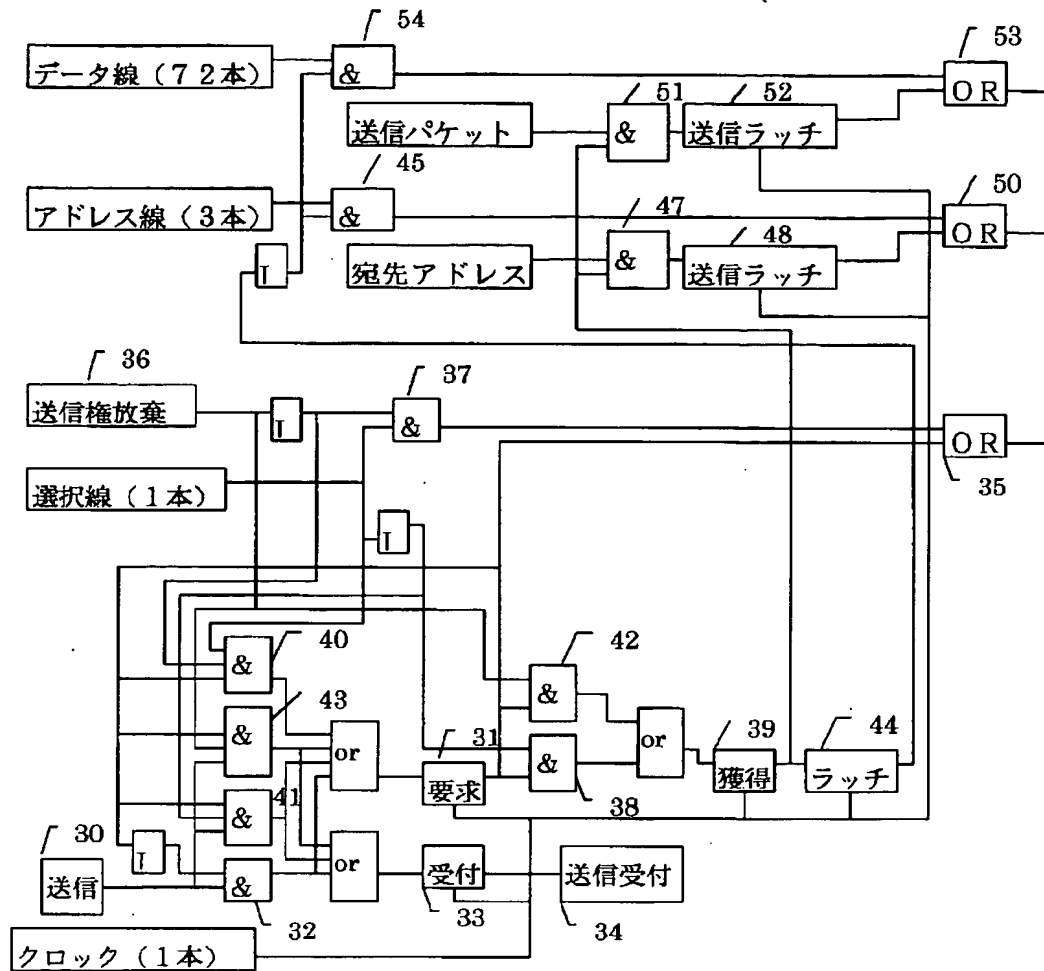
図 4 0





【図 38】

図 38



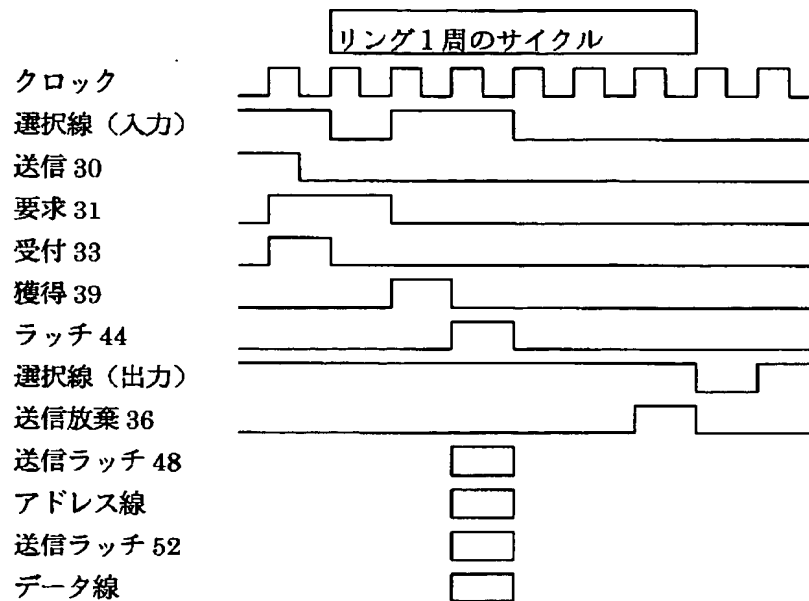
【図 41】

図 41



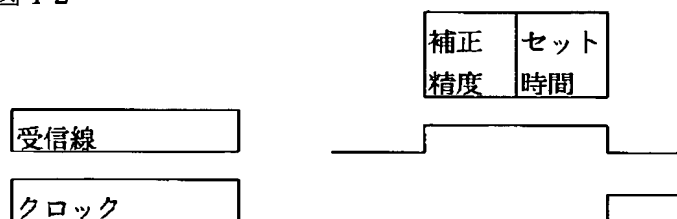
【図 3 9】

図 3 9



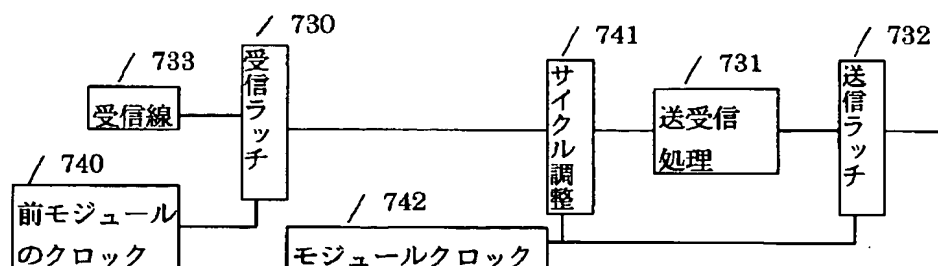
【図 4 2】

図 4 2



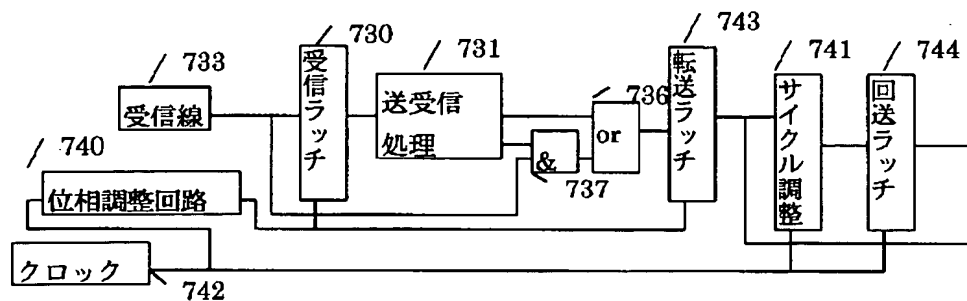
【図 4 3】

図 4 3



【図 4 4】

図 4 4



【図 4 5】

【図 6 2】

図 4 5

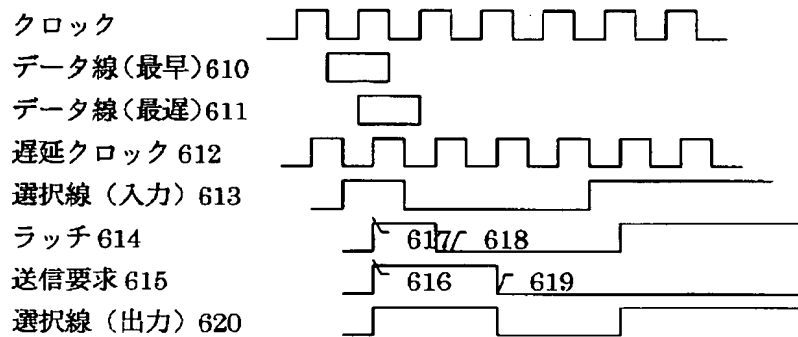
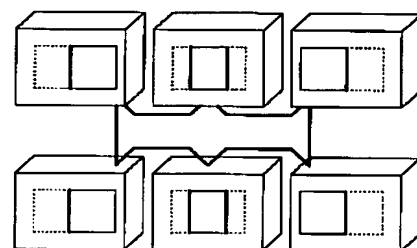
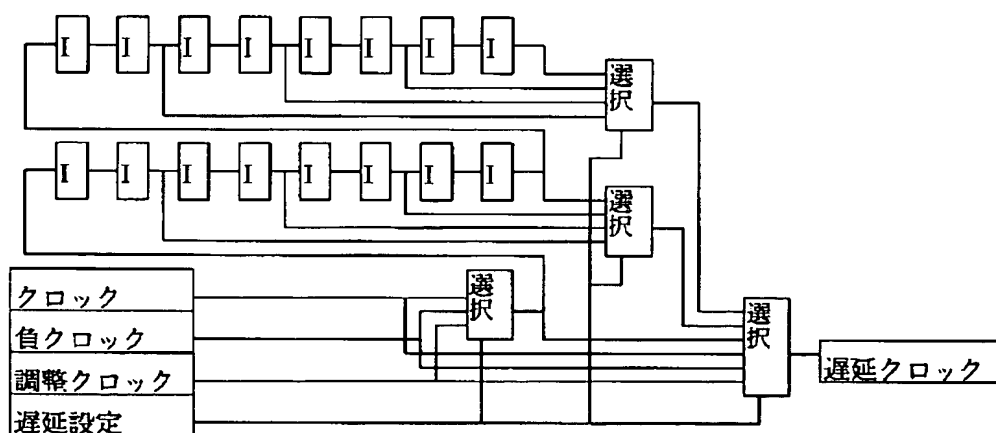


図 6 2

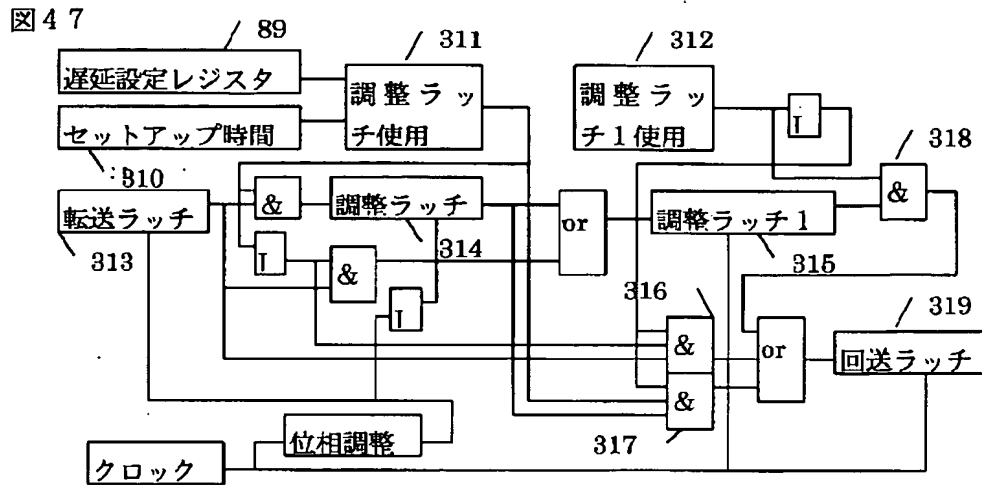


【図 4 6】

図 4 6



【図 47】

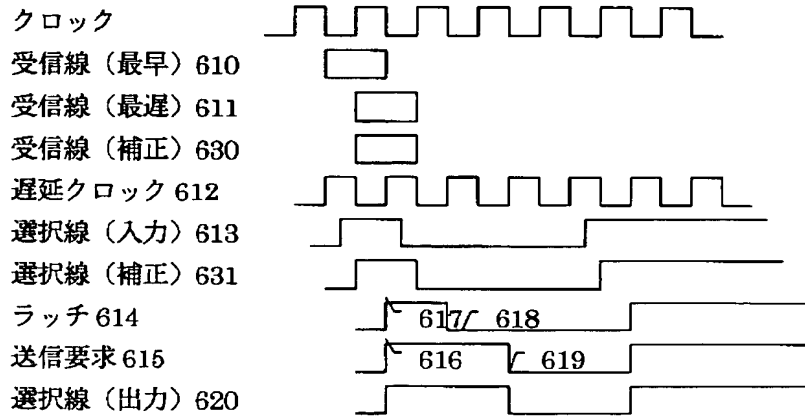


【図 48】

【図 69】

図 48

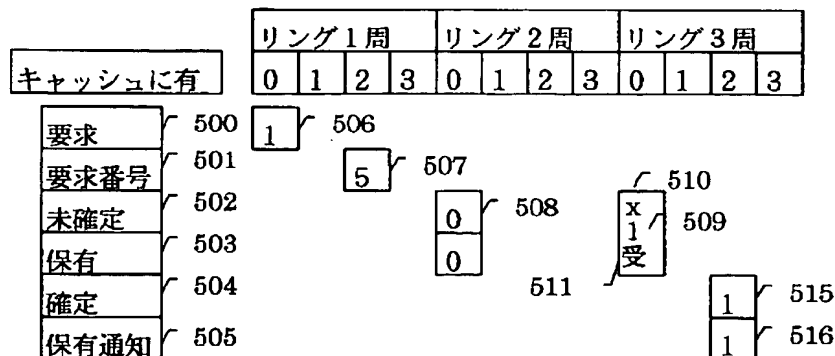
図 69



写し無し (N)
写し有り (C)
写しベクトル有り (CH)
変更有り (M)
写しベクトル有り (MH)

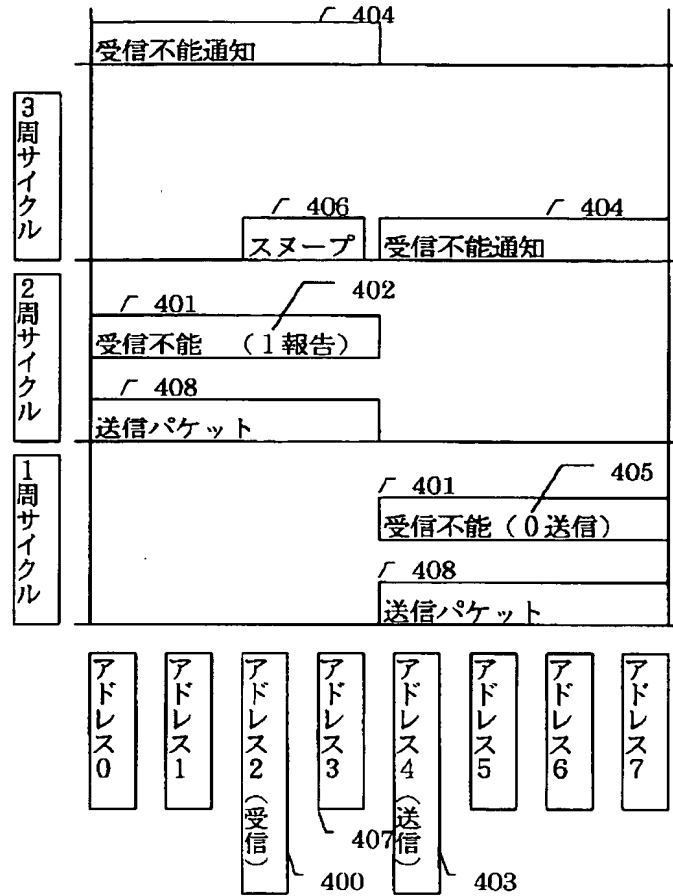
【図 52】

図 52



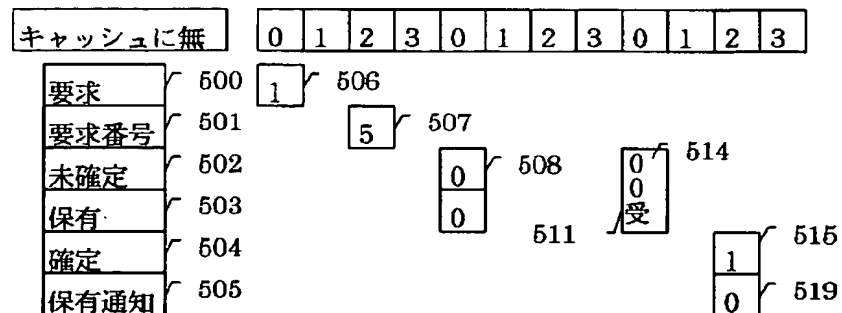
【図 4 9】

図 4 9



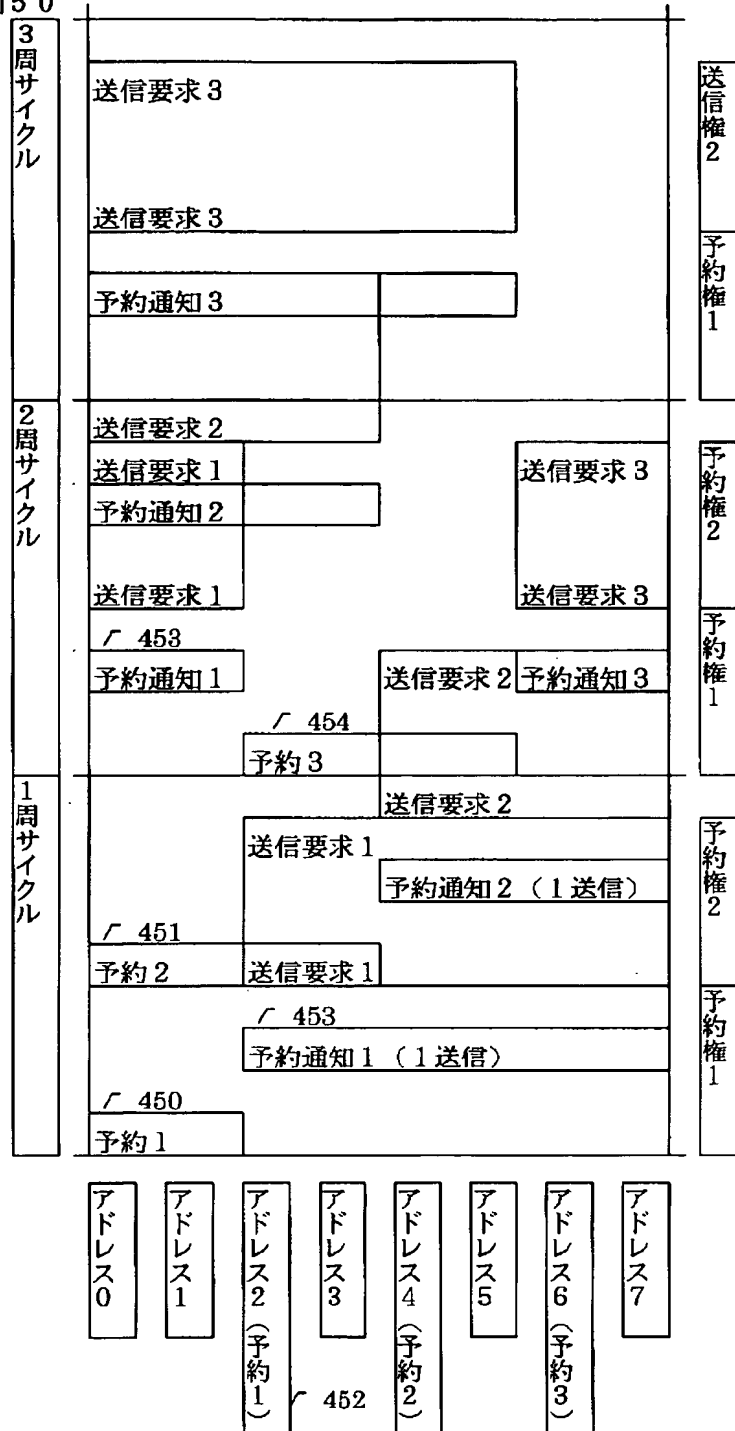
【図 5 3】

図 5 3



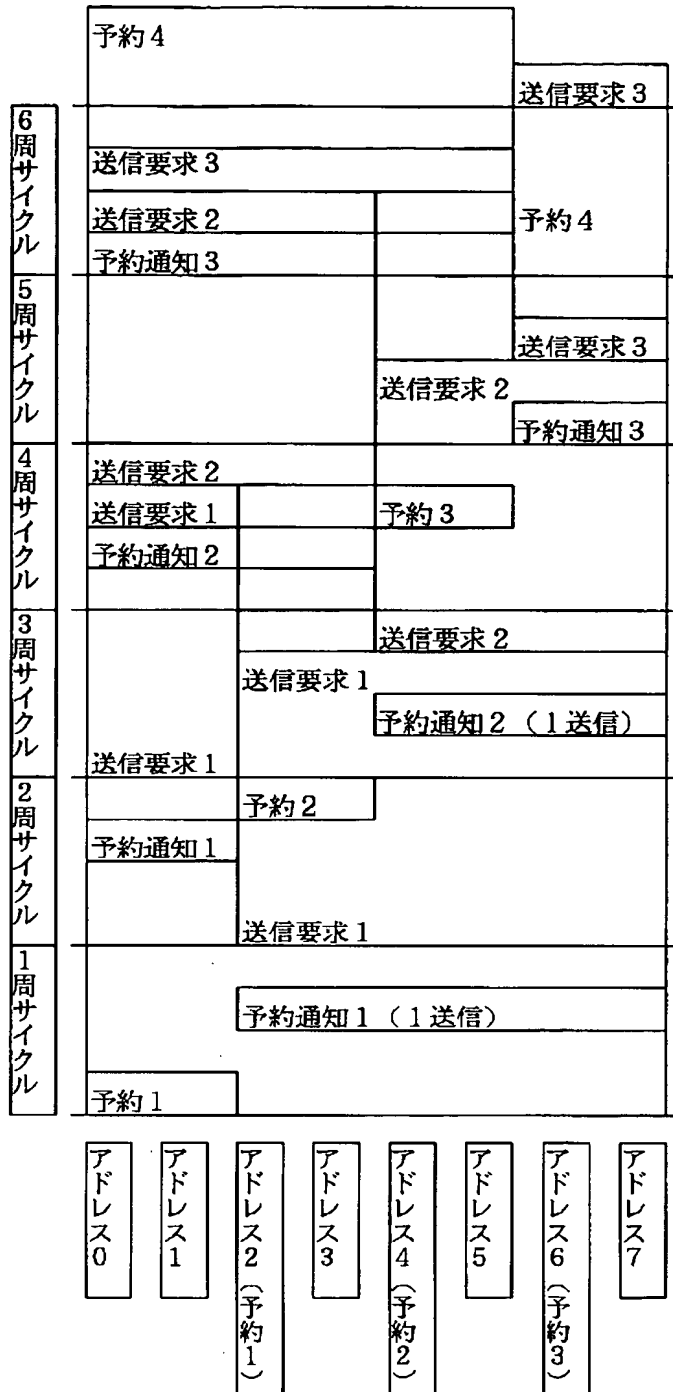
【図 5 0】

図 5 0



【図 5 1】

図 5 1



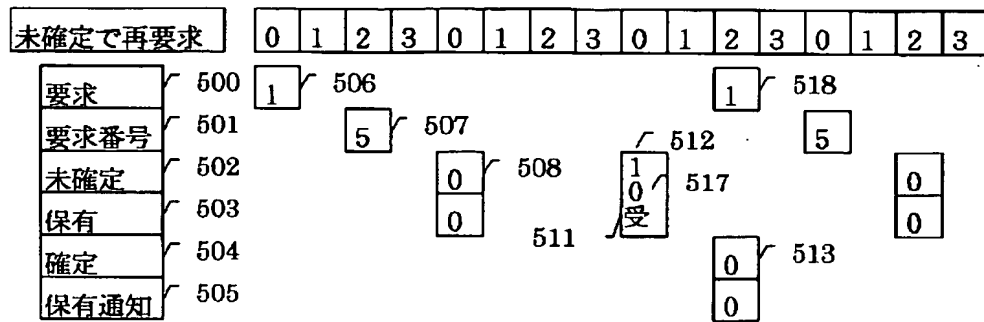
【図 6 8】

図 6 8

1	アドレス (4 / 8 ビット)	写しベクトル (8 ビット)
2	アドレス (4 / 8 ビット)	写しベクトル (8 ビット)
4M	アドレス (4 / 8 ビット)	写しベクトル (8 ビット)

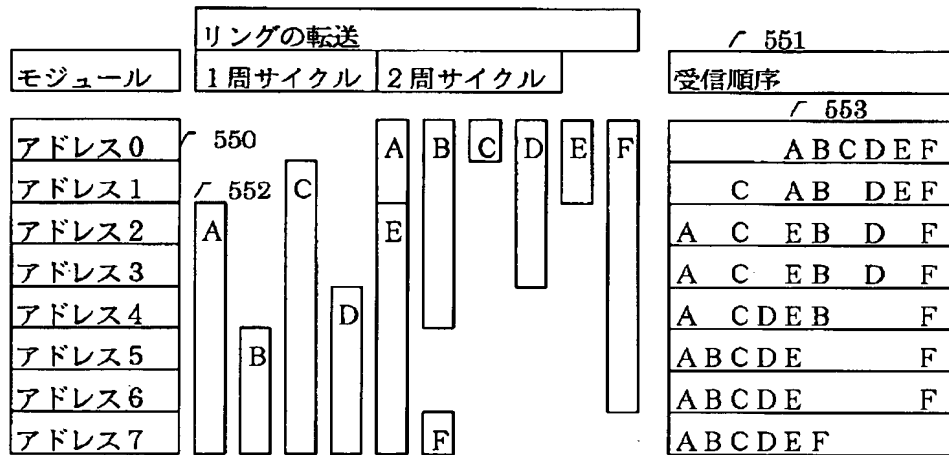
【図 5 4】

図 5 4



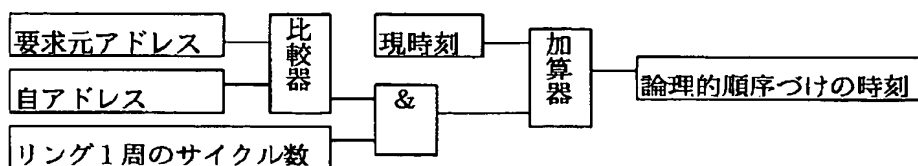
【図 5 5】

図 5 5



【図 5 6】

図 5 6





【図 5 7】

【図 6 6】

図 5 7

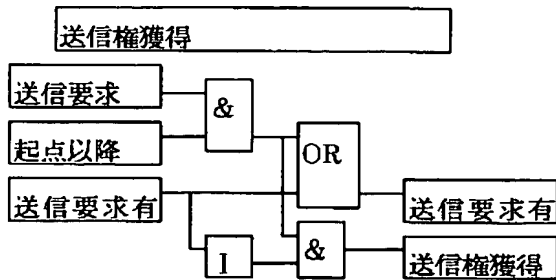
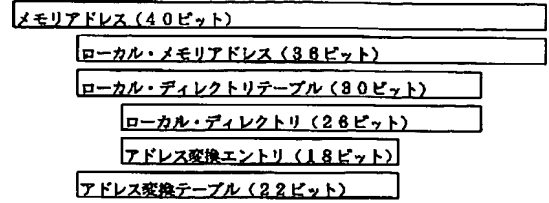
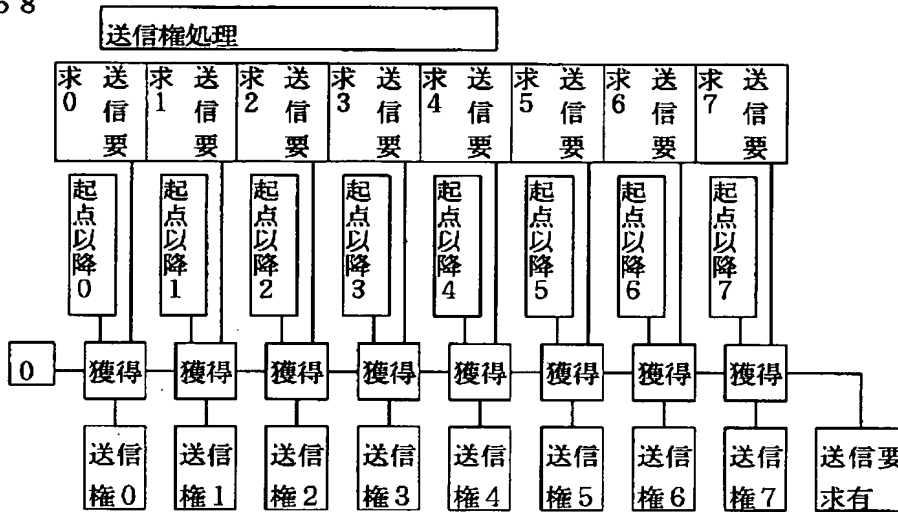


図 6 6



【図 5 8】

図 5 8



【図 6 0】

図 6 0

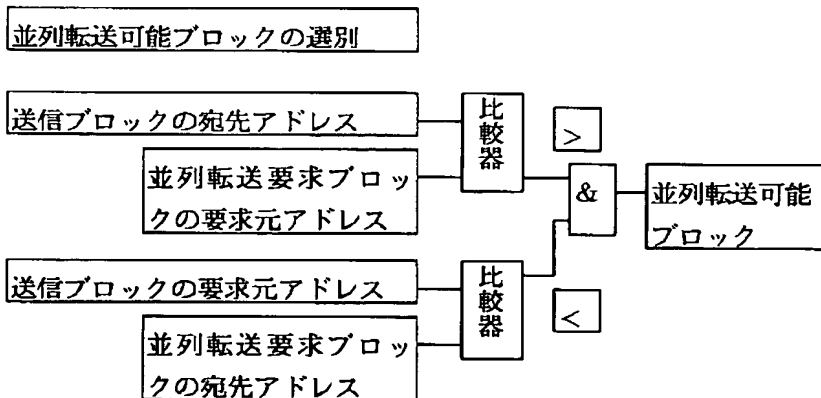


图 5 9

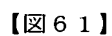
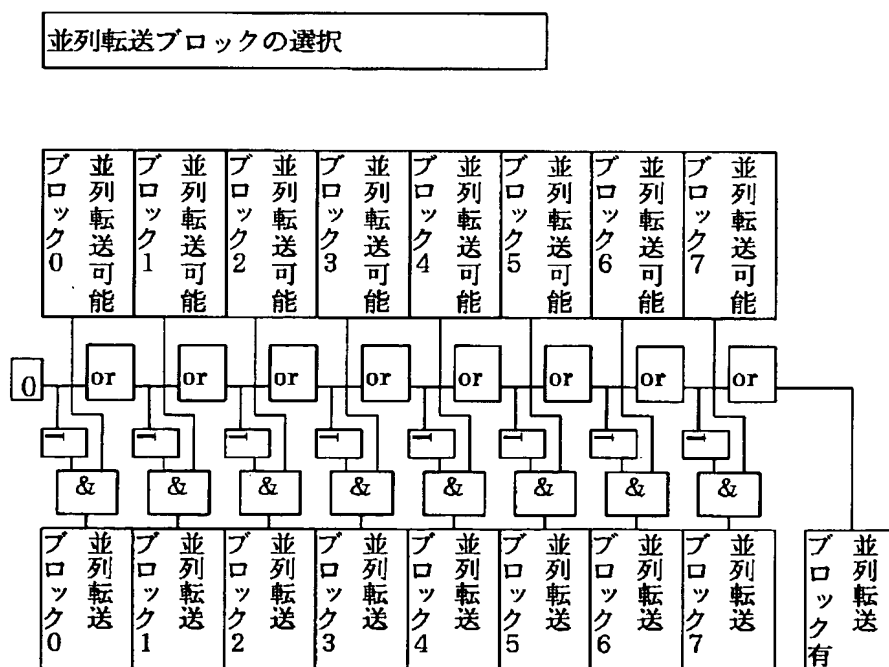
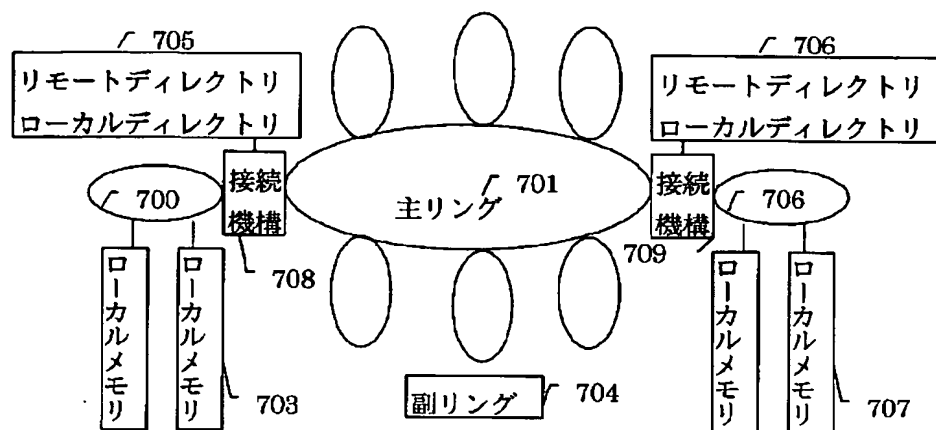


图 6-1



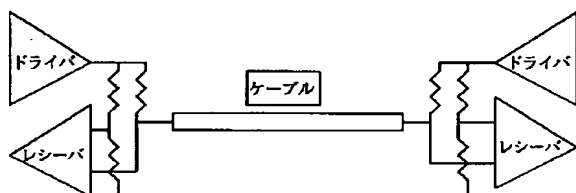
【図 6 3】

図 6 3



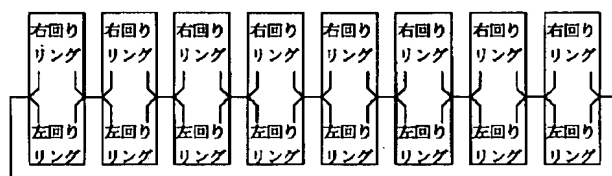
【図 7 0】

図 7 0



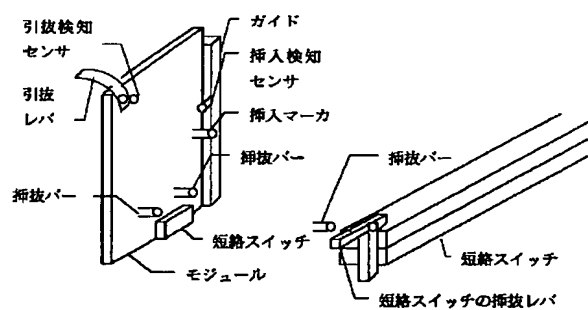
【図 7 1】

図 7 1



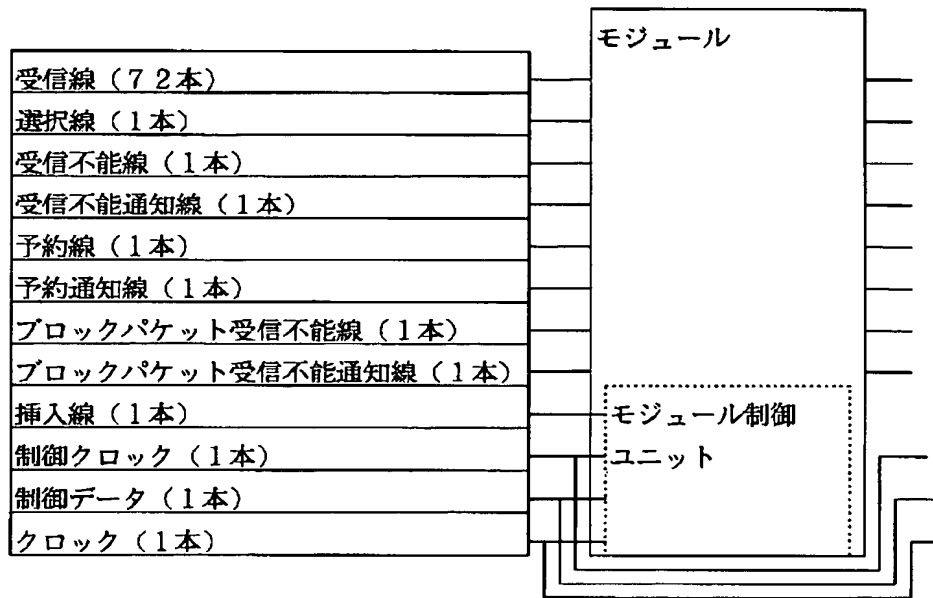
【図 7 2】

図 7 2



【図 7 3】

図 7 3



**This Page is Inserted by IFW Indexing and Scanning  
Operations and is not part of the Official Record**

## **BEST AVAILABLE IMAGES**

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ BLACK BORDERS
- ☐ IMAGE CUT OFF AT TOP, BOTTOM OR SIDES
- ☐ FADED TEXT OR DRAWING
- ☒ BLURRED OR ILLEGIBLE TEXT OR DRAWING
- ☒ SKEWED/SLANTED IMAGES
- ☐ COLOR OR BLACK AND WHITE PHOTOGRAPHS
- ☐ GRAY SCALE DOCUMENTS
- ☐ LINES OR MARKS ON ORIGINAL DOCUMENT
- ☐ REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY
- ☐ OTHER: \_\_\_\_\_

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.**